Institute of Education



CLS Missing Data Strategy

Richard Silverwood & Michalis Katsoulis

27 April 2023

CENTRE FOR LONGITUDINAL STUDIES



Economic and Social Research Council

Outline



- 1. CLS Missing Data Strategy
- 2. 1958 National Child Development Study (NCDS)
- 3. 1970 British Cohort Study (BCS70)
- 4. Other CLS studies
- 5. COVID-19 Surveys
- 6. Leveraging data linkages
- 7. Resources

CLS Missing Data Strategy

CLS Missing Data Strategy



- All users of longitudinal data need to consider the issue of missing data since some non-response is inevitable.
- Strategies for how to deal with missing data depend on the nature of non-response.
- Well known (principled) methods for handling missing data include multiple imputation, inverse probability weighting and full information maximum likelihood.
- These rely on the assumption that the data are missing at random (MAR), implying that systematic differences between the missing values and the observed values can be explained by observed data.

CLS Missing Data Strategy



- Most studies employing such MAR methods rely on a largely arbitrary selection of variables used as predictors of missingness.
- We aim to maximise the plausibility of the MAR assumption by optimising the set of such variables used in analyses.
- We use systematic (data driven) approaches to identify variables that are associated with non-response at each sweep in each study.
- This allows us to capitalise on the rich data cohort members have provided over the years/decades in order to deal with missing data and reduce bias.

Why the focus on (wave) non-response?

- CENTRE FOR LONGITUDINAL STUDIES
- Wave non-response is the main driver of missing data in analyses of CLS studies. Item non-response less of an issue.
- Much of the wave non-response is due to attrition.
- For longitudinal analyses, wave non-response at the most recent sweep is therefore usually the biggest contributor to missingness.
- Can identify predictors of wave non-response at cohort (rather than analysis) level pragmatic approach.
- In analyses in which item non-response is more prevalent, this may need additional consideration.



1958 National Child Development Study (NCDS)







Identifying predictors of non-response in NCDS



- Aim to maximise the plausibility of the MAR assumption by exploiting the richness of NCDS data.
- Using a data driven approach we identify the variables that are associated with non-response at each sweep.
- These can then be used as auxiliary variables.
- Substantive interest in understanding the drivers of non-response within and between cohorts.

Identifying predictors of non-response in NCDS

- ~17,500 variables in NCDS sweeps 0-8.
- Exclude:
 - Routed variables.
 - Binary variables with prevalence <1%.
 - Variables with item non-response > 50%.
- Use summary scores for scales.
- Use summary measures; exclude constituent variables.
- 587 variables meeting inclusion criteria
 - \rightarrow multi-stage, data driven approach.

Identifying predictors of non-response in NCDS

- For non-response at sweep *t*:
 - Stage 1: Univariable regressions for predictors at sweep 0, ..., sweep t – 1. Retain predictors with p < 0.05.
 - Stage 2: Multivariable regressions for predictors at sweep 0, ..., sweep t 1. Retain predictors with p < 0.05.
 - Stage 3: MI. Multivariable regressions for predictors at sweep 0, ..., sweep t – 1, adjusted for predictors at previous waves. Retain predictors with p < 0.001.
- Full details in Mostafa et al (2021) and NCDS Missing Data User Guide.

Predictors of non-response



	NR sweep 1 (age 7)	NR sweep 2 (age 11)	NR sweep 3 (age 16)	NR sweep 4 (age 23)	NR sweep 5 (age 33)	NR sweep 6 (age 42)	NR BM sweep (age 44)	NR sweep 7 (age 46)	NR sweep 8 (age 50)	NR sweep 9 (age 55)
Sweep 0 (birth)	3	1	1	4	3	3	5	3	3	6
Sweep 1 (age 7)		5	3	3	5	1	5	4	3	4
Sweep 2 (age 11)			1	4	3	3	1	3	2	2
Sweep 3 (age 16)				4	4	3	4	4	4	5
Sweep 4 (age 23)					5	2	1	2	3	2
Sweep 5 (age 33)						5	4	2	3	5
Sweep 6 (age 42)							5	3	5	2
BM sweep (age 44)								3	3	1
Sweep 7 (age 46)									1	1
Sweep 8 (age 50)										3
Total	3	6	5	15	20	17	25	24	27	31

Consistent predictors of participation



- Being female (adulthood only).
- Lower childhood social class in childhood; higher childhood social class in adulthood.
- Higher early life cognitive ability; fewer adolescent conduct problems.
- Social participation; voting; union membership.
- Being married; home ownership.
- Participation in previous sweeps.

Simple test analyses

We test the performance of the missing data strategy using two approaches:

- 1. "Travelling back in time" to see whether distributions of variables from earlier sweeps can be replicated using only data from respondents at a later sweep.
- 2. Comparison to external population benchmarks.







Social class of mother's husband at birth







Multiple imputation

- Complete case analysis only (generally) valid under MCAR.
- Many simple imputation approaches are problematic:
 - Mean imputation.
 - Last observation carried forward.
 - Single conditional imputation.
- In MI, plausible values are used in place of the missing values in a way that allows:
 - 1. Parameter estimates to be unbiased.
 - 2. Uncertainty to be estimated in an appropriate way.
- MI valid under MAR.

Multiple imputation



- Specify an appropriate imputation model.
- Create a series of imputed datasets.
- Each imputed dataset analysed using the substantive model.
- Results combined using standard rules.
- Can be undertaken using standard statistical software.
- Widely adopted as practical for applied researchers in a wide range of settings.

Which variables should be included in the imputation model?

Definitely:

• All variables in the substantive model, including any interactions.

Optional "auxiliary variables":

- Variables associated with the underlying values of the variable(s) subject to missingness.
- Particularly those also associated with the probability of missingness.









_





























CENTRE FOR

No educational qualifications at age 50



35

CENTRE FOR

STUDIES

LONGITUDINAL






More realistic analyses

- Such simple analyses are useful for testing/illustrating the basic idea.
- More realistic analyses likely to be more complicated:
 - More variables in substantive model exposure(s), outcome(s), control variables,...
 - Inclusion of auxiliary variable(s) predictive of non-response at further sweeps.
 - Inclusion of auxiliary variable(s) predictive of the underlying values of the variable(s) subject to missingness.
 - Inclusion of auxiliary variable(s) to deal with item nonresponse.
 - Different types of substantive model.

More realistic analyses



- Basic idea remains the same.
- Main concern likely to be (in the MI setting) instability of the imputation model caused by number/type of variables.
- Illustrative realistic example used throughout the NCDS Missing Data User Guide (next session).



- We have identified variables which predict non-response at each sweep of NCDS.
- These can be used as auxiliary variables in subsequent analyses to increase the plausibility of the MAR assumption.
- Simple test analyses have shown this approach to perform well.
- A straightforward approach, easily implemented in standard software.
- Lists of predictors of non-response available via NCDS Missing Data User Guide.
- Will be updated when new sweeps of data become available.
- Also work using linked data (see later).

1958 National Child Development Study (NCDS)

CENTRE FOR LONGITUDINAL STUDIES



Abstract

Objective: Non-response is unavoidable in longitudinal surveys. The consequences are lower statistical power and the potential for bias. We implemented a systematic data-driven approach to identify predictors of non-response in the National Child Development Study (NCDS; 1958 British birth cohort). Such variables can help make the missing at random assumption more plausible, which has implications for the handling of missing data

Study Design and Setting: We identified predictors of non-response using data from the 11 sweeps (birth to age 55) of the NCDS (n = 17,415), employing parametric regressions and the LASSO for variable selection.

Results: Disadvantaged socio-economic background in childhood, worse mental health and lower cognitive ability in early life, and lack of civic and social participation in adulthood were consistently associated with non-response. Using this information, along with other data from NCDS, we were able to replicate the "population distribution" of educational attainment and marital status (derived from external data), and the original distributions of key early life characteristics.

Conclusion: The identified predictors of non-response have the potential to improve the plausibility of the missing at random assumption. They can be straightforwardly used as "auxiliary variables" in analyses with principled methods to reduce bias due to missing data. © 2021 Elsevier Inc. All rights reserved.

Keywords: Cohort studies; Longitudinal data; Missing data; Multiple imputation; National Child Development Study; Non-response

1. Introduction

Non-response is unavoidable in longitudinal surveys. The consequences are smaller samples due to attrition, lower statistical power and decreased representativeness compared to the originally intended target population. With some exceptions where complete case analysis is valid [1-3], in the majority of analyses of longitudinal data bias will occur if the implications of selection due to incompleteness are not formally addressed [4,5]. There is a broad interdisciplinary consensus that missing data should be dealt with using principled approaches and it has recently been argued that "complete-case analysis should be used with the same caution we ascribe to unadjusted estimates, as its validity relies on strong, often unrealistic assumptions" [6].

Rubin described 3 missing data generating mechanisms: i) missing completely at random (MCAR); ii) missing

Conflict of interest: None. * Corresponding author. E-mail address: R Silverwood@ucl.ac.uk (R I. Silverwood)

https://doi.org/10.1016/j.jclinepi.2021.02.019 0895-4356/© 2021 Elsevier Inc. All rights reserved. at random (MAR); iii) missing not at random (MNAR) [3,7,8]. MCAR implies that the probability of non-response does not depend on any variable (measured or unmeasured), or that there are no systematic differences between the observed and missing data. MCAR is partially testable, since we can examine whether variables available in our data are associated with missingness. MAR implies that systematic differences between the missing values and the observed values can be explained by observed data, or that given the observed data, the reasons for missingness do not depend on unobserved variables. With some exceptions for specific missing data patterns [9,10] the MAR assumption is untestable [11]. The third mechanism - MNAR - implies that that the observed data are insufficient to explain variation in the probability of missingness. MNAR is also untestable and methods to deal with this type of missing data generating mechanism rely heavily on further - usually distributional - assumptions [12]

Contextualizing the 1958 British National Child Development Study (NCDS) within Rubin's framework, we know that the missing data generating mechanism is not Institute of Education



Handling missing data in the

National Child Development

Study

User guide (Version 2)

July 2021





CENTRE FOR LONGITUDINAL STUDIES

1970 British Cohort Study (BCS70)

Identifying predictors of non-response in BCS



- Very similar approach used in BCS as in NCDS
- Aim to maximise the plausibility of the MAR assumption using a data driven approach we identify the variables that are associated with non-response at each sweep (and can potentially be used as auxiliary variables)
- We also highlight cases in which we can explore MNAR

Identifying predictors of non-response in NCDS

- ~20000 variables in BCS sweeps 0-8.
- Exclude:
 - Routed variables.
 - Binary variables with prevalence <1%.
 - Variables with item non-response > 40%.
- Use summary scores for scales
- For non-response at sweep *t* we used the same 3 stage approach as in NCDS (using a bit stricter criteria)

Predictors of non-response



	NR sweep 1 (age 7)	NR sweep 2 (age 11)	NR sweep 3 (age 16)	NR sweep 4 (age 23)	NR sweep 5 (age 33)	NR sweep 6 (age 42)	NR sweep 7 (age 44)	NR sweep 8 (age 46)	NR sweep 9 (age 50)
Sweep 0 (birth)	4	3	1	5	1	5	4	2	5
Sweep 1 (age 7)		4	1	3	2	2	3	1	0
Sweep 2 (age 11)			1	3	1	1	3	1	1
Sweep 3 (age 16)				0	1	0	0	0	0
Sweep 4 (age 23)					2	3	2	1	0
Sweep 5 (age 33)						2	0	2	2
Sweep 6 (age 42)							0	0	2
Sweep 7 (age 46)								1	2
Sweep 8 (age 46)									1
Total	4	7	3	11	7	13	12	8	13

Consistent predictors of participation

- Being female (adulthood only).
- Few household moves
- Paternal social class (early sweeps)
- Higher early life cognitive ability;
- Social participation voting;
- Home ownership.
- Participation in previous sweeps.

CENTRE FOR

UDIFS



Internal validation: Cognitive ability - at 7yo



CENTRE FOR LONGITUDINAL STUDIES

Internal validation: Cognitive ability - at 7yo



CENTRE FOR

STUDIES

LONGITUDINAL

Mean BMI levels – Age 34 (MEN)





SENSITIVITY ANALYSIS – External validation BMI levels – Age 34 (MEN)



CENTRE FOR

STUDIES

LONGITUDINAL

SENSITIVITY ANALYSIS – External validation: Mean BMI levels



We followed the same procedure for

- Women at age 34
- Men at age 42
- Women at age 42



- We have identified variables which predict non-response at each sweep of BCS.
- These can be used as auxiliary variables in subsequent analyses to increase the plausibility of the MAR assumption.
- This approach can be extended for MNAR mechanisms in some cases, with appropriate external benchmark
- Simple test analyses have shown this approach to perform well.
- A straightforward approach, easily implemented in standard software.

CENTRE FOR LONGITUDINAL STUDIES

Other CLS studies

Next Steps



Institute of Education A data driven approach to understanding and handling non-response in the Next Steps cohort CLS vorking paper number 20205

By Richard J. Silverwood, Lisa Calderwood, Joseph W Sakshaug, George B. Ploubidis





Millennium Cohort Study (MCS)



_

• Ongoing work.

CENTRE FOR LONGITUDINAL STUDIES

COVID-19 Surveys

Non-response weights



- To correct for non-response in the COVID-19 surveys, non-response weights are provided, so that IPW analyses can be undertaken.
- Non-response weights capitalise on the rich data cohort members have provided over many years.





	Wave 2							
Cohort	Issued sample	Response within issued sample	Target population	Response within target population				
NSHD	2551	1569 (61.5%)	3758	1488 (39.6%)				
NCDS	11,655	6282 (53.9%)	15,291	6228 (40.7%)				
BCS70	12,133	5320 (43.9%)	17,486	5236 (29.9%)				
Next Steps	11,529	3664 (31.8%)	15,770	3609 (22.9%)				
MCS cohort members	13,547	3274 (24.2%)	19,243	3233 (16.8%)				
Total	51,415	20,109 (39.1%)	71,548	19,794 (27.7%)				

Target population: individuals born in the specified birth period who are alive and still residing in the UK.

Overall response rate within issued sample (39.1%) comparable to similar COVID-19 web surveys.

Derivation of non-response weights



- 1. Within sample corresponding to target population, model COVID-19 survey response conditional on a common set of covariates using logistic regression.
- 2. For COVID-19 survey respondents, predict probability of response from model.
- 3. Calculate non-response weight as inverse of probability of response.
- 4. Examine distribution of weights across cohorts to decide whether truncation may be desirable; apply truncation if so.
- 5. Calibrate weights so they sum to number of respondents in each cohort.

Derivation of non-response weights Response model



- Selection of covariates in response model informed by literature and results of the CLS Missing Data Strategy, plus assumed associations with the probability of response and/or with key COVID-19 survey variables.
- Aimed to use broadly same set of variables in each cohort to ensure consistency.
- Not possible to include identical sets of variables due to data being collected at different ages and using different questions.
- Full details in COVID-19 Surveys User Guide.

Derivation of non-response weights Response model



Sex

Ethnicity Parental social class Number of rooms at

home/persons per room

Cognitive ability

Early life mental health

Voting

Membership in organisations

Internet access prior to web survey Consent for biomarkers Consent for linkages Educational qualifications Economic activity Partnership status Psychological distress BMI

Self-rated health **Smoking status** Maternal mental health Social capital/social support Income Number of nonresponses across all previous sweeps **Response at COVID-19** Wave 1 survey*



Grey: all cohort members; red: COVID-19 Wave 2 survey respondents only; blue: COVID-19 Wave 2 survey respondents after application of non-response weights.

Comparison with MI (work in progress)



Grey: all cohort members; red: COVID-19 Wave 2 survey respondents only; blue: after application of non-response weights; green: after application of non-response weights.

CENTRE FOR

STUDIES

LONGITUDINAL

CENTRE FOR LONGITUDINAL STUDIES

Leveraging data linkages

Extending the strategy

- Growing interest in whether linked administrative data have the potential to aid analyses subject to missing data in cohort studies.
- Identify predictors of cohort non-response in linked administrative data.
- Explore whether added value in including identified variables as auxiliary variables with respect to restoring sample representativeness.
- Focusing on linked NCDS and hospital episode statistics (HES) data here. Many other linkages with CLS cohort data available.

Hospital Episode Statistics (HES)



- A collection of databases containing details of interactions with NHS hospitals in England.
- Linkage between NCDS and HES datasets undertaken on the basis of consent at sweep 8 (age 50).
- Matching conducted using deterministic linkage based on combinations of the participant's name, sex, date of birth and postcode.
- Linked data available via secure access through the UK Data Service.

HES predictors of NCDS non-response



- A total of 58 variables derived from HES data relating to:
 - Numbers of admissions and appointments
 - Missed appointments
 - Investigations undertaken
 - Diagnoses
 - Treatments received
- Employed a similar approach to identify most important predictors of NCDS non-response at wave 9 (age 55).
- 10 variables identified.

Restoring NCDS sample representativeness

- Undertook similar test analysis to see if including the identified HES variables helped restore sample representative.
- Concluded that it did a bit but essentially no additional gain relative to using only previously identified survey predictors of non-response.

FNTRF FOR

NGITUDINAL

Leveraging data linkages





Statistics data to aid the handling of non-response and restore sample representativeness in the **1958 National Child Development** Study

CLS working paper number 2023/1

Nasir Rajah¹, Lisa Calderwood¹, Bianca L De Stavola², Katie Harron², George B Ploubidis¹ and Richard J Silverwood^{1*}

1. Centre for Longitudinal Studies, UCL Social Research Institute, 20 Bedford Way, London WC1H 0AL 2. Population, Policy & Practice Research and Teaching Department, UCL Great Ormond Street Institute of Child Health, 30 Guilford Street, London WC1N 1EH





Research Council

CENTRE FOR LONGITUDINAL STUDIES

Resources
Resources



 Handling missing data webpage: <u>https://cls.ucl.ac.uk/data-access-training/handling-missing-data/</u>



Resources

- Mostafa T, Narayanan M, Pongiglione B, Dodgeon B, Goodman A, Silverwood RJ, Ploubidis GB. <u>Missing at random assumption</u> <u>made more plausible: evidence from the 1958 British birth</u> <u>cohort</u>. J Clin Epidemiol. 2021;136:44-54.
- Silverwood RJ, Goodman A, Ploubidis GB. Letter to the editor: Don't forget survey data: 'healthy cohorts' are 'real-world' relevant if missing data are handled appropriately. Longitudinal and Life Course Studies. 2022;13(2):335-41.
- Silverwood R, Narayanan M, Dodgeon B, Ploubidis G. <u>Handling</u> <u>missing data in the National Child Development Study: User</u> <u>Guide (Version 2)</u>. London: UCL Centre for Longitudinal Studies; 2021.

Resources



- Silverwood RJ, Calderwood L, Sakshaug JW, Ploubidis GB. <u>A data</u> <u>driven approach to understanding and handling non-response in the</u> <u>Next Steps cohort</u>. CLS Working Paper 2020/5. London: UCL Centre for Longitudinal Studies; 2020.
- Brown M, Goodman A, Peters A, Ploubidis GB, Sanchez A, Silverwood R, et al. <u>COVID-19 Survey in Five National Longitudinal</u> <u>Studies: Waves 1, 2 and 3 User Guide (Version 3)</u>. London: UCL Centre for Longitudinal Studies and MRC Unit for Lifelong Health and Ageing; 2021.
- Rajah N, Calderwood L, De Stavola BL, Harron K, Ploubidis GB, Silverwood RJ. <u>Using linked Hospital Episode Statistics data to better</u> <u>handle non-response and restore sample representativeness in the</u> <u>National Child Development Study</u>. CLS Working Paper Series 2023/1. London: UCL Centre for Longitudinal Studies; 2023.

Institute of Education



Thank you.

CENTRE FOR LONGITUDINAL STUDIES