# Handling missing data in the 1970 British Cohort Study

Richard Silverwood, Michalis Katsoulis & Brian Dodgeon

6 June 2024

CENTRE FOR LONGITUDINAL STUDIES

UK RI
Economic and Social Research Council

CENTRE FOR
LONGITUDINAL
STUDIES

Introduction to BCS70

**1946**

**19 46** National Survey of Health and Development

NHS

Photo: BiblioArchives/ LibraryArchives

ncds
National Child Development Study

**1958**

**1970**

BCS**70**
1970 British Cohort Study

NEXT STEPS
LEARNING FROM YOUR GENERATION

**1989**

DIG DEEP FOR THE MINERS
NATIONAL UNION OF MINEWORKERS

**2000**

MILLENNIUM COHORT STUDY

There is NO PLANET B

I VOTED LEAVE

Photo: Miguelon756-5303

Now Playing
6 of 157
Keep in the Dark
Temples
Sun Structures
2:24

http://www.

EARLY LIFE COHORT FEASIBILITY STUDY

**2022**

# 1970 British Cohort Study (BCS70)

- Longitudinal birth cohort study of people born in England, Wales, Scotland and Northern Ireland in one week in 1970.

- Initial N = 17,196, though NI members not followed up after birth (unless moved to mainland).

- Boost samples through childhood.

- Collected information on health, physical, educational and social development, and economic circumstances, among other factors.

- PI: George Ploubidis.

# BCS70: Respondents, instruments and response

| | 1970 Birth | 1975 5 | 1980 10 | 1986 16 | 1996 26 | 2000 30 | 2004 34 | 2008 38 | 2012 42 | 2016 46 |
|---|---|---|---|---|---|---|---|---|---|---|
| main respondent | mother | parents | subject/ parents | subject/ parents | subject | subject | subject | subject | subject | subject |
| secondary respondent | medical | medical | medical/ school | medical/ school | | | children | | | medical |
| survey instruments | | cognitive assess-ments | cognitive assess-ments | cognitive assess-ments and diaries | 16-page postal questionnaire | | cognitive assess-ments | telephone survey | vocabulary test | cognitive assess-ments |
| linked data | | | | | | | | | Consent to data linkage | |
| response | 17,196 | 13,135 | 14,875 | 11,622 | 9,003 | 11,261 | 9,665 | 8,874 | 9,841 | 8,581 |

# BCS70: Topics covered by life stage

| Birth | School years | Adult |
|---|---|---|
| Family | Family | Family (partners, children) |
| Parental employment | Parental employment | Employment |
| Obstetric history | Financial circumstances | Income |
| Smoking in pregnancy | Housing | Housing |
| Pregnancy (problems, antenatal care) | Health | Health |
| Labour (length, pain relief, problems) | Behaviour | Health-related behaviour |
| Birth (problems, sex, weight, length) | School | Courses and qualifications |
| | Views and expectations | Basic skills |
| | Attainment | Cognitive ability |
| | | Views and expectations |

# Age 51 (available November 2024)

## Employment and income

- Occupation
- Income
- Partner's employment and income
- Impact of COVID-19
- Benefits
- Pensions
- Debts
- Intergenerational transfers

## Physical health and health behaviours

- General
- Longstanding illness
- Health conditions
- COVID tests, symptoms, long COVID
- Height/weight
- Exercise
- Diet
- Drinking and smoking

## Cognitive skills and processes

- Immediate and delayed recall
- Animal naming
- Letter cancellation
- National Audit Reading Test (NART)

## Mental health and well-being

- Psychological distress
- Mental well-being
- Life satisfaction

## Family and relationships

- Grandchildren
- CM's and partners parents – caring responsibilities
- Social contact
- Quality of relationships
- Menstruation
- Fertility

## Activities, attitudes and values

# Cohort profiles

Sullivan A, Brown M, Hamer M, Ploubidis G. Cohort Profile Update: The 1970 British Cohort Study (BCS70). International Journal of Epidemiology. 2023; 52(3): e179-e86.

Elliott J, Shepherd P. Cohort Profile: 1970 British Birth Cohort (BCS70). International Journal of Epidemiology. 2006; 35(4): 836-43.

# Accessing more information/data



https://cls.ucl.ac.uk



https://ukdataservice.ac.uk/

CENTRE FOR
LONGITUDINAL
STUDIES

CLS missing data strategy

# Missing data

- Non-response inevitable in longitudinal surveys.

- Reduced analysis sample size → reduced efficiency.

- Respondents often systematically different from non-respondents → threat to representativeness & potential for bias.

- 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) = given the observed values, missingness does not depend on unobserved values.

# Multiple imputation (MI)

- Specify an appropriate imputation model and create a series of imputed datasets.

- Each imputed dataset is analysed using the substantive model and results combined using standard rules.

- Need to include all variables in the substantive model in the imputation model.

- Can also include "auxiliary variables": variables associated with the underlying values of the variable(s) subject to missingness, particularly those also associated with the probability of missingness.

# CLS Missing Data Strategy

- Most analyses 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 wave.

- These can then be considered for inclusion as auxiliary variables.

- This allows us to capitalise on the rich data cohort members have provided over the years/decades.

# Why the focus on wave non-response?

- 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

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. 1958 British birth cohort. Journal of Clinical Epidemiology 2022; 136: 44-54.

# Next Steps

## RESEARCH ARTICLE

### A data-driven approach to understanding non-response and restoring sample representativeness in the UK Next Steps cohort

Richard J. Silverwood, R.Silverwood@ucl.ac.uk
Lisa Calderwood, L.Calderwood@ucl.ac.uk
Morag Henderson, Morag.Henderson@ucl.ac.uk
University College London, UK

Joseph W. Sakshaug, joe.sakshaug@iab.de
University of Warwick, UK
Institute for Employment Research, Germany
Ludwig Maximilian University of Munich, Germany

George B. Ploubidis, G.Ploubidis@ucl.ac.uk
University College London, UK

Non-response is common in longitudinal surveys, reducing efficiency and introducing the potential for bias. Principled methods, such as multiple imputation, are generally required to obtain unbiased estimates in surveys subject to missingness which is not completely at random. The inclusion of predictors of non-response in such methods, for example as auxiliary variables in multiple imputation, can help improve the plausibility of the missing at random assumption underlying these methods and hence reduce bias. We present a systematic data-driven approach used to identify predictors of non-response at Wave 8 (age 25–26) of Next Steps, a UK national cohort study that follows a sample of 15,770 young people from age 13–14 years. The identified predictors of non-response were across a number of broad categories, including personal characteristics, schooling and behaviour in school, activities and behaviour outside of school, mental health and well-being, socio-economic status, and practicalities around contact and survey completion. We found that including these predictors of non-response as auxiliary variables in multiple imputation analyses allowed us to restore sample representativeness in several different settings, though we acknowledge that this is unlikely to universally be the case. We propose that these variables are considered for inclusion in future analyses using principled methods to explore and attempt to reduce bias due to non-response in Next Steps. Our data-driven approach to this issue could also be used as a model for investigations in other longitudinal studies.

**Keywords** cohort studies • missing data • multiple imputation • non-response • sample representativeness

Silverwood RJ, Calderwood L, Henderson M, Sakshaug JW, Ploubidis GB. **A data-driven approach to understanding non-response and restoring sample representativeness in the UK Next Steps cohort**. Longitudinal and Life Course Studies. 2024; 15(2): 227-50.

17

# BCS70

A data driven approach to address missing data in the 1970 British birth cohort

Michail Katsoulis[1], Martina Narayanan[2], Brian Dodgeon[2], George Ploubidis[2], Richard Silverwood[2]

1. MRC unit for Lifelong Health and Ageing, Institute of Cardiovascular Science, UCL
2. Centre for Longitudinal Studies, Institute of Education, UCL

Katsoulis M, Narayanan M, Dodgeon B, Ploubidis G, Silverwood R. **A data driven approach to address missing data in the 1970 British birth cohort**. medRxiv. 2024: 2024.02.01.24302101.

# Extending the strategy

Rajah N, Calderwood L, De Stavola BL, Harron K, Ploubidis GB, Silverwood RJ. **Using linked administrative data to aid the handling of non-response and restore sample representativeness in cohort studies: the 1958 national child development study and hospital episode statistics data**. BMC Medical Research Methodology. 2023; 23(1): 266.

# Additional resources

Missing data strategy in BCS70

# Identifying predictors of non-response in BCS70

- Very similar approach used in BCS70 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)

# Summary

- We have identified variables which predict non-response at each wave of BCS70.

- 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 Handling missing data in the CLS cohort studies (https://cls.ucl.ac.uk/wp-content/uploads/2020/04/Handling-Missing-Data-User-Guide-2024.pdf).

- Will be updated when new waves of data become available.

# BCS: Response over time

- Individuals who were not resident in the UK in 1970 (here included as not respondents) were added later between age 5 to age 16 (sweep 4) and hence affecting response rates between sweeps 1 to 3

# Identifying predictors of non-response in BCS70

- ~20000 variables in BCS70 sweeps 1-10

- 18037 individuals in BCS70

- 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)

# 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 $p < 0.001$.

  - Stage 2: Multivariable regressions for predictors at sweep 0, …, sweep $t - 1$ (separately). 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 $p < 0.001$

Variables sex, country of birth, participation in all previous sweeps and father's socioeconomic status are contained in the final set of variables and hence are not included in this formula

# Flowchart presenting the number of variables used to find the predictors of non-response

**Number of available variables: N=21021**

Excluded "routed" variables (questions that depend on a specific response to a previous question), variables with >40% missing, variables with categories including >98% of cases, recoded variables with categories <1%, excluded all binary variables with prevalence <1% of cases if possible, and used index/score variables that combined information from many variables rather the individual constituent items

**Variables to be used for Stage 1-3: N=967**

Application of the Stage 1-3 process

**Variables to be used after Stage 3:**

**Vary from N=7 (at sweep 2) to**

**N=16 (at sweep 10)**

# Predictors of non-response

| | NR sweep 2 (age 5) | NR sweep 3 (age 10) | NR sweep 4 (age 16) | NR sweep 5 (age 26) | NR sweep 6 (age 30) | NR sweep 7 (age 34) | NR sweep 8 (age 38) | NR sweep 9 (age 42) | NR sweep 10 (age 46) |
|---|---|---|---|---|---|---|---|---|---|
| Sweep 1 (birth) | S1: 28 S2: 15 | S1: 26 S2: 14 | S1: 30 S2: 13 | S1: 37 S2: 15 | S1: 30 S2: 14 | S1: 34 S2: 13 | S1: 35 S2: 13 | S1: 31 S2: 12 | S1: 35 S2: 11 |
| Sweep 2 (age 5) | | S1: 15 S2: 8 | S1: 39 S2: 9 | S1: 56 S2: 6 | S1: 49 S2: 9 | S1: 55 S2: 11 | S1: 57 S2: 7 | S1: 49 S2: 8 | S1: 51 S2: 9 |
| Sweep 3 (age 10) | | | S1: 64 S2: 15 | S1: 88 S2: 14 | S1: 84 S2: 12 | S1: 91 S2: 13 | S1: 90 S2: 13 | S1: 82 S2: 9 | S1: 78 S2: 9 |
| Sweep 4 (age 16) | | | | S1: 63 S2: 11 | S1: 39 S2: 7 | S1: 57 S2: 9 | S1: 65 S2: 12 | S1: 49 S2: 9 | S1: 54 S2: 7 |
| Sweep 5 (age 26) | | | | | S1: 35 S2: 7 | S1: 54 S2: 8 | S1: 64 S2: 13 | S1: 49 S2: 5 | S1: 53 S2: 10 |
| Sweep 6 (age 30) | | | | | | S1: 56 S2: 11 | S1: 75 S2: 11 | S1: 56 S2: 12 | S1: 64 S2: 15 |
| Sweep 7 (age 34) | | | | | | | S1: 80 S2: 10 | S1: 69 S2: 17 | S1: 77 S2: 16 |
| Sweep 8 (age 38) | | | | | | | | S1: 23 S2: 10 | S1: 27 S2: 8 |
| Sweep 9 (age 42) | | | | | | | | | S1: 39 S2: 5 |
| STAGE 3 without extra variables* | S3: 5 (out of 15 in total from S2) | S3: 7 (out of 22 in total from S2) | S3: 11 (out of 37 in total from S2) | S3: 10 (out of 46 in total from S2) | S3: 9 (out of 49 in total from S2) | S3: 12 (out of 65 in total from S2) | S3: 12 (out of 79 in total from S2) | S3: 9 (out of 82 in total from S2) | S3: 13 (out of 90 in total from S2) |
| Total with extra variables† | 7 | 9 | 14 | 13 | 12 | 15 | 15 | 14 | 16 |

*We do not include in these counts the variables sex, country of birth, participation in all previous sweeps and father's socioeconomic status which were used directly in Stage 3

†We include in these counts the variables sex, country of birth, participation in all previous sweeps and father's socioeconomic status

28

## Consistent predictors of participation

- Method of contraception (early sweeps)

- Paternal completion of education (early sweeps)

- Higher early life cognitive ability (sweeps in adulthood)

- Being female (sweeps in adulthood only)

- Few household moves (sweeps in adulthood only)

- Social participation – voting (sweeps in adulthood only)

- Home ownership. (sweeps in adulthood only)

- Parity, i.e. number of older siblings

- Participation in previous sweeps

- Paternal social class

# Internal validation: Paternal Social class - at birth

# Internal validation: Paternal Social class - at birth

# Internal validation: Paternal Social class - at birth

- Chained equations.

- 50 imputations.

- Auxiliary variables: All predictors of non-response at sweep 10 (age 46) from sweeps 0-9.

32

# Internal validation: Cognitive ability - at 5yo

# Internal validation: Cognitive ability - at 5yo

# Internal validation: Cognitive ability - at 5yo

- Chained equations.

- 50 imputations.

- Auxiliary variables: All predictors of non-response at sweep 10 (age 46) from sweeps 1-9.

35

# Hazard ratios of father's social class on participants' mortality after age 26 i) in the original sample ii) among 8314 respondents (complete case) of sweep 5 (age 26) and iii) after MI

# Hazard ratios of father's social class on participants' mortality after age 26 i) in the original sample ii) among 8314 respondents (complete case) of sweep 5 (age 26) and iii) after MI

# Hazard ratios of father's social class on participants' mortality after age 26 i) in the original sample ii) among 8314 respondents (complete case) of sweep 5 (age 26) and iii) after MI

# Tackling missing data in BCS70
**IN PRACTICE**

To address problems due to missing data in BCS70

- Utilise information from https://cls.ucl.ac.uk/wp-content/uploads/2020/04/Handling-Missing-Data-User-Guide-2024.pdf, pages 88-110
- Based on the sweep you want to tackle missing data, select the predictors of non-response you need
- E.g. when using MI, these variables can be used as auxiliary variables in your imputation model

# Tackling missing data in BCS70
## IN PRACTICE – sweep 2 (age 5)

**Table B1.** Predictors of non-response at Sweep 2 (age 5).

| Sweep | Variable description | Variable | Variable derivation details |
|---|---|---|---|
| Sweep 1 (age 0) | Marital Status | a0012 | Single<br>Married<br>Divorced/Separated |
| | Parity | A0166_new | Recoded from A0166<br><br>0<br>1<br>2<br>3<br>>4 |
| | Father's social status | BD1BPOS_new | Recoded from BD1BPOS<br><br>I Single, no work, unskilled other<br>II partial work<br>III manual work<br>III non-manual work<br>IV managerial/technical work<br>V professorial work |
| | Country of Birth | COB_new | Recoded from COB<br><br>England<br>Wales<br>Scotland |

| Sweep | Variable description | Variable | Variable derivation details |
|---|---|---|---|
| | | | Other |
| | Father's age at completion of education | A0010_new | Recoded from a0010<br><br>≤15 year old<br>16-18 year old<br>≥19 years old |
| | Number of antenatal visits | a0190 | Continuous (per visit) |
| | Method of contraception | a0029b | None<br>Pill alone<br>Pill alone and other method<br>Other method |

# Tackling missing data in BCS70
## IN PRACTICE – sweep 5 (age 26)

**Table B4**. Predictors of non-response at Sweep 5 (age 26).

| Sweep | Variable description | Variable | Variable derivation details |
|---|---|---|---|
| Sweep 1 (age 0) | Father's social status | BD1BPOS_new | Recoded from BD1BPOS<br><br>I Single, no work, unskilled other<br>II partial work<br>III manual work<br>III non-manual work<br>IV managerial/technical work<br>V professorial work |
| | Age of mother at 1st birth | BD1AGEFB | Continuous (per year) |
| | Method of contraception | a0029b | None<br>Pill alone<br>Pill & Other method<br>Other Method |
| | Parity (i.e. number of older siblings) | A0166_new | Recoded from a0166<br><br>0<br>1<br>2<br>3<br>>4 |
| | Certainty of last menstrual period | a0196 | Certain vs uncertain |

| Sweep | Variable description | Variable | Variable derivation details |
|---|---|---|---|
| | Sex | SEX | Female vs male |
| Sweep 2 (age 5) | External score | Extern_score_5 | Per unit increase |
| | Harris scoring method | f114 | Per unit increase |
| Sweep 3 (age 10) | Gross family income | grfaminc | Per unit increase |
| | Teacher Rutter assessment | B3T_Rutt | Per unit increase |
| | Number of household accessories | b3hldstf | Per unit increase |
| Sweep 4 (age 16) | Satisfaction with teen's school progress | Pb3_1 | Very satisfied<br>Fairly satisfied<br>Not satisfied |
| Non-response at Previous sweeps (i.e. sweeps 1, 2 & 3) | | | Yes vs no |

# Tackling missing data in BCS70

## IN PRACTICE – sweep 10 (age 46)

**Table B9**. Predictors of non-response at Sweep 10 (age 46).

| Sweep | Variable description | Variable | Variable derivation details |
|---|---|---|---|
| Sweep 1 (age 0) | Father's social status | BD1PSOC_new | Recoded from BD1PSOC<br><br>I Single, no work, unskilled other<br>II partial work<br>III manual work<br>III non-manual work<br>IV managerial/technical work<br>V professorial work |
| | Country of Birth | COB_new | Recoded from COB<br><br>England<br>Wales<br>Scotland<br>Other |
| | Parity (i.e. number of older siblings) | A0166_new | Recoded from a0166<br><br>0<br>1<br>2<br>3<br>>4 |
| | Certainty of last menstrual period | a0196 | Certain vs uncertain |
| | Was lactation attempted | a0297 | Not attempted vs attempted |

# Tackling missing data in BCS70
## IN PRACTICE – sweep 10 (age 46) - continued

| Sweep | Variable description | Variable | Variable derivation details |
|---|---|---|---|
| | Number of antenatal visits | a0190 | Per visit |
| Sweep 2 (age 5) | Copying designs score | f119 | Continuous (Per unit) |
| Sweep 3 (age 10) | Score BAS Matrices | BASmatrx | Continuous (Per unit) |
| | Accommodation | D2_new | Recoded from d2<br><br>Owned<br><br>Bought<br><br>Council rented<br><br>Other rented<br><br>Tied to occupation |
| Sweep 4 (age 16) | | | |
| Sweep 5 (age 26) | | | |
| Sweep 6 (age 29) | Had eczema or skin problems? | Othskin_new | Recoded from othskin (8,9--<br>>missing)<br><br>No vs yes |
| | Voted in general elections 1997? | Vote97_new | Recoded from vote97 (8, 9 --<br>>missing)<br><br>No vs yes |
| Sweep 7 (age 34) | Is this address participant's residence? | b7nrmal | No vs yes |
| Sweep 8 (age 38) | Willing to be contacted for Parents Research Project | b8parent | No vs Yes |
| | Any children aged 0-6 | b8chd006 | No vs Yes |

# Tackling missing data in BCS70

## IN PRACTICE – sweep 10 (age 46) - continued

| Sweep | Variable description | Variable | Variable derivation details |
|---|---|---|---|
| Sweep 9 (age 42) | Total score | B9VSCORE | Per mark |
| Non-response at Previous sweeps (i.e. sweeps 1-8) | | | Yes vs no |

These variables can be used to restore representativeness and

- Present summary statistics of the characteristics of individuals

- Assess causal relationships between your exposure and your outcome, account for potential confounders

# Summary

- We have identified variables which predict non-response at each wave of BCS70.

- 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.

- Lists of predictors of non-response available via <span style="color:red">Handling missing data in the CLS cohort studies</span> (https://cls.ucl.ac.uk/wp-content/uploads/2020/04/Handling-Missing-Data-User-Guide-2024.pdf).

- Will be updated when new waves of data become available.

- Our work will facilitate researchers who plan to utilise BCS70 and help them address bias due to missing data

Q&A

Thank you

**UCL**

CENTRE FOR
LONGITUDINAL
STUDIES

Economic
and Social
Research Council