

# Longitudinal Data and Research – A Deep Dive

Vernon Gayle  
Professor of Sociology & Social Statistics  
University of Edinburgh

vernon.gayle@ed.ac.uk  
@profbigvern  
2024

# The Research Value of Longitudinal Data

Vernon Gayle  
Professor of Sociology & Social Statistics  
University of Edinburgh

vernon.gayle@ed.ac.uk  
@profbigvern  
2024

© Vernon Gayle

# Why Use Longitudinal Data?

- UK has an unparalleled collection
- These resources are critical for analysing social change (and social stability)
- But they need justification because they are costly in money and time

# Longitudinal Social Surveys

- Cross-sectional data
  - Respondents surveyed at only one time point
- Longitudinal data
  - Repeated contacts (with the same individuals)
  - Respondents surveyed at multiple time points



# Longitudinal Social Science Study Designs

## Panel Study

The panel are the group and are repeatedly studied

- US (PSID)
- Germany (SOEP)
- Britain BHPS/UKHLS
- Australia (HILDA)
- Canada (SLID)
- Swiss (SHP); Korea (KLIPS); Russia (RLMS)

# Longitudinal Social Science Study Designs

## Cohort Study

- Repeated contacts data collection  
(simply a specific form of panel design in my view)
- Principally concerned with charting the development  
of a particular ‘group’ from a certain point in time

# Longitudinal Social Science Study Designs

- Cohort Study

- A birth cohort of babies born in a particular year (e.g. 1946; 1958; 1970; 2000-2)
- A youth cohort, a group of pupils who completed compulsory education in the same year (YCS; LSYPE)

# Research Using Longitudinal Social Survey Datasets

- For many social research projects cross-sectional data will be sufficient
- Most social research projects can be improved by the analysis of longitudinal data
- Some research questions require longitudinal data

# Questions that Require Longitudinal Data

- Flows into and out of poverty
- The effects of family migration on the woman's subsequent employment activities
- Numerous policy intervention examples
- Numerous examples relating to 'individual' development

# Key Messages (so far....)

- For many social research projects cross-sectional data will be sufficient
- Most social research projects can be improved by the analysis of longitudinal data
- *Researchers are likely to make more rapid progress using existing large-scale longitudinal data resources*

- Some research questions require longitudinal data
- Longitudinal data are not a panacea



*'This longitudinal study suggests that notwithstanding the dominant effect of severity of intellectual impairment, a number of factors within and outside the family may also contribute to higher attainment in reading, writing and numeracy.'*

*In particular mainstream schooling for those with less severe disabilities appears to have benefited the children in this study' (p.390).*

Turner, S., Alborz, A. and **Gayle, V.** (2008) 'Predictors of academic attainments of young people with Down's syndrome', *Journal of Intellectual Disability Research*, 52(5), pp. 380-392.

# Subjective Well-Being & Happiness

- Non-economic measures of social progress
- “Improving the quality of our lives should be the ultimate target of public policies” Angel Gurría, OECD Secretary-General
- UK commitment to developing wider measures of well-being
- Tailoring government policies to the things that matter



- Moving house itself causes a boost in happiness, and brings people back to their initial levels
- Moving and set-point theory
- Long-distance migrants are at least as happy as short-distance migrants despite the higher social and psychological costs involved
- Re-theorize moving within a conceptual framework that accounts for social well-being from a life-course perspective

Nowok, B., van Ham, M., Findlay, A. and **Gayle, V.** (2013) 'Does migration make you happy? A longitudinal study of internal migration and subjective wellbeing', *Environment and Planning A*, 45(4), pp. 986-1002.

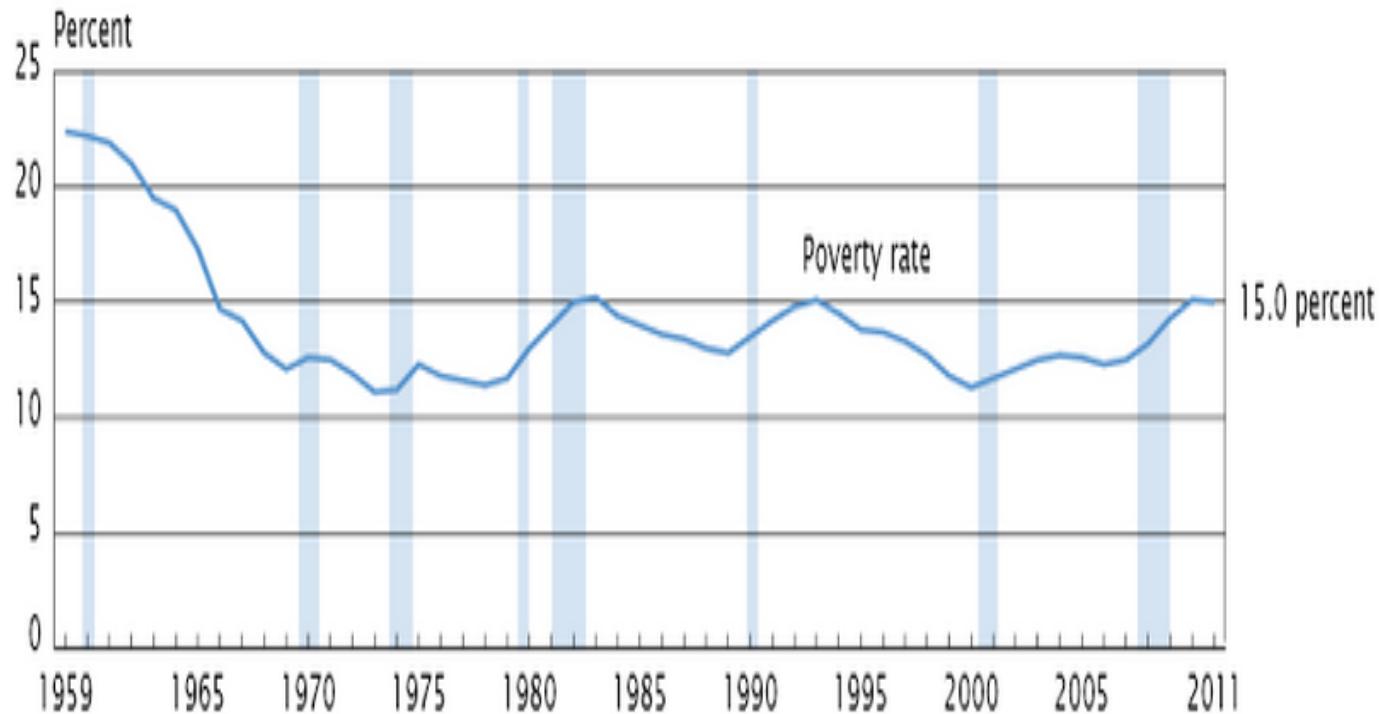
# The Bigger Picture

- UKHLS is the largest living observatory of contemporary social life
- Contribution to the ‘evidence base’
- Contribution to empirically informed planning
- *Influencing behaviour and informing interventions*
- *Contributing to a fair and vibrant society*

# Examples...

- Cohort studies - secondary smoking effects on children
- Whitehall Studies - influenced successive governments' thinking on social gradients in health
- Whitehall Studies -dispelled the myth that high status jobs have higher risk of heart disease

# USA Poverty Rate 1959 - 2011



Note: The data points are placed at the midpoints of the respective years. For information on recessions, see Appendix A.

Source: U.S. Census Bureau, Current Population Survey, 1960 to 2012 Annual Social and Economic Supplements.

- Poverty rates flattened out in 1990s
- BHPS showed apparent cross-sectional stability but a hidden longitudinal flux
  - Substantial turnover or churning
  - The poor were not always poor
- Not detectable without panel data!

- UK Poverty rate approximately 18%
- In a 6 year periods one-third of individuals were poor at least once
- Only 2% were poor for all six years!
- Repeated short spells of poverty were more common than one long spell

# The Consequences...

- Contributed to the ‘rubber band theory’
  - we are attached to an elastic tether
- Influenced the Labour government’s welfare reforms in the late 1990s
  - focussing on moving people into work and making work pay
- Now influences how living standards are measured in Britain
  - Official Statistics now include household panel based information

# Summary Messages

- For many social research projects cross-sectional data will be sufficient
- Most social research projects can be improved by the analysis of longitudinal data
- Some research questions require longitudinal data

# A vignette...

The story of Jason Jones (aged 10) and his mum



# Questions that Require Longitudinal Data

- Flows into and out of poverty
- The effects of family migration on the woman's subsequent employment activities
- Numerous policy intervention examples
- Numerous examples relating to 'individual' development

# Methodological Benefits of Longitudinal Social Science Data

- Micro-level social processes
- Temporal ordering of events
- Improving control for residual heterogeneity
- Improving control for state dependence

# Micro-Level Social Processes

- Cross-sectional data = a snap shot
  - Good for studying the immediate
  - Several datasets can study macro / or gross changes
- Repeated contacts data allow the study of
  - The passage of time
  - Individual (or household) change/stability
  - Processes that occur at the micro-level of the individual (or family)
  - Surprises (or shocks)

# Temporal Ordering of Events (Direction of Influence)

- Time moves in one direction so...
  - An event in 1990 comes before an event in 1995
  - Experiences at primary school could affect university entry
  - Teenage smoking could influence health in old age
- But not *vice versa*  
*One sociology professor has argued with me suggesting that time does not move in only one direction*

# Temporal Ordering of Events (Direction of Influence)

- There is unequivocal evidence from cross-sectional data that, overall, the unemployed have poorer health
- This is consistent with both
  - A. Unemployment causing ill health
  - B. Ill health causing unemployment
- These two substantive stories are quite different

<b>Month</b>	<b>Level of Health (20 = Good Health)</b>	<b>Employ Status</b>
1	17	Employed
2	17	Employed
3	17	Employed
4	17	Unemployed
5	17	Unemployed
6	10	Unemployed
7	16	Unemployed
8	5	Unemployed
9	4	Unemployed
10	3	Unemployed
11	2	Unemployed
12	1	Unemployed

# Person A



Became unemployed this has affected  
his level of health

<b>Month</b>	<b>Level of Health (20 = Good Health)</b>	<b>Employ Status</b>
1	17	Employed
2	1	Employed
3	1	Employed
4	1	Unemployed
5	1	Unemployed
6	1	Unemployed
7	1	Unemployed
8	1	Unemployed
9	1	Unemployed
10	1	Unemployed
11	1	Unemployed
12	1	Unemployed

# Person B



Poor health led to unemployment  
(because of poor job performance)

# In a cross-sectional study (at month 12)

- Person A would have been unemployed for 9 months and have a health score of 1
- Person B would have been unemployed for 9 months and have a health score of 1
- This is an obvious example of how panel (i.e. repeated contacts) data can make an essential contribution to untangling social relationships

# Improving Control for Omitted Explanatory Variables

- Residual Heterogeneity
  - Omitted explanatory variables
  - Unobserved heterogeneity
- The possibility of substantial variation between similar individuals due to unmeasured, and possibly immeasurable, variables is known as '*residual heterogeneity*'

# Improving Control for Omitted Explanatory Variables

*Because data collection instruments often fail to capture the detailed nature of social life there is, almost inevitably, considerable heterogeneity in response variables even amongst respondents that share the same characteristics across all of the explanatory variables*

# Improving Control for Omitted Explanatory Variables

*As long as we make the assumption that (at least some of) these effects are enduring there are techniques for accounting for omitted explanatory variables if we have data at more than one time point*

# Improving Control for Omitted Explanatory Variables

- There are no routine methods of accounting for omitted explanatory variables in cross-sectional analysis
- It is sometimes claimed that the main advantage of longitudinal data is that it facilitates improved control for the plethora of variables that are omitted from any analysis
- Panel data won't completely sweep this problem away, but suitable models can improve control for, and estimate the effects of, residual heterogeneity

# Improving Control for the Effects of Previous States (state dependence)

*A frequently noted empirical regularity in the analysis of unemployment data is that those who were unemployed in the past or have worked in the past are more likely to be unemployed (or working) in the future*

(Nobel Prize winner J.J. Heckman)

# Improving Control for the Effects of Previous States (state dependence)

- Much of human behaviour is influenced by previous behaviour and outcomes (positive feedback)
- McGinnis (1968) '*axiom of cumulative inertia*'

# Improving Control for the Effects of Previous States (state dependence)

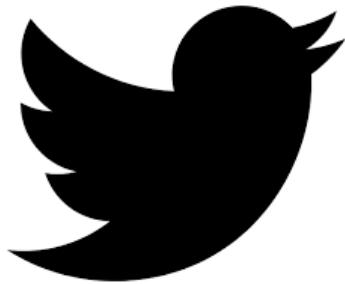
- Working in May = more likely to be working in June
- Married this year = more likely to be married next year
- Own your own house this quarter
- Travel to work by car this week

# Improving Control for the Effects of Previous States (state dependence)

With panel data we may be able to include past behaviour in the modelling process

# Summary Message

There are methodological benefits...  
but panel data are not a panacea!



Tweet - Longitudinal data enhance  
our ability to investigate complicated  
processes in the social world

# Sources of Longitudinal Data

Vernon Gayle  
Professor of Sociology & Social Statistics  
University of Edinburgh

vernon.gayle@ed.ac.uk  
@profbigvern  
2024

© Vernon Gayle

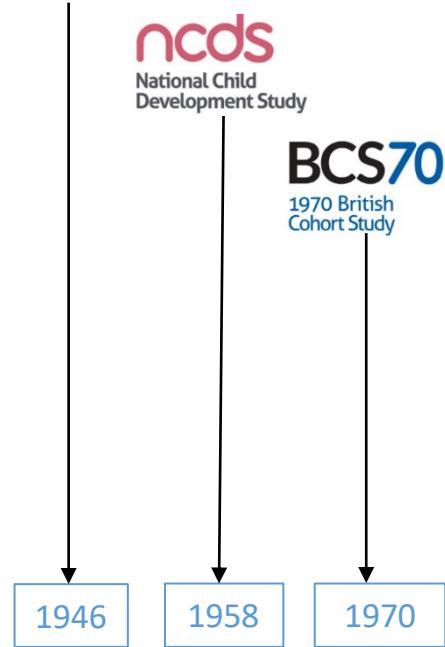
# Studying Longer Term Trends

- Many sources of ‘repeated’ cross-sectional data
  - Rapid progress can be made
  - Standard statistical approached (e.g. regression models)
- 
- Comparability (equivalence) is the central challenge
  - How should time be represented

# Cohort Studies

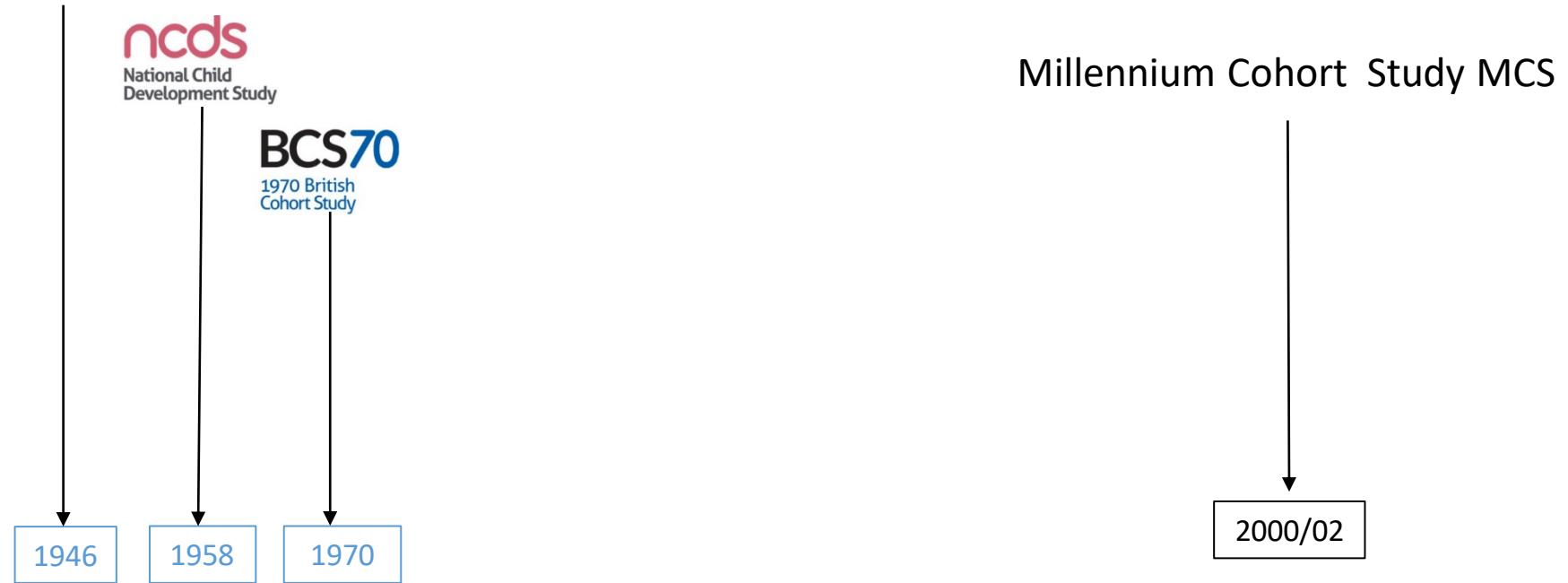
# UK Birth Cohorts

National Survey of Health  
and Development 1946



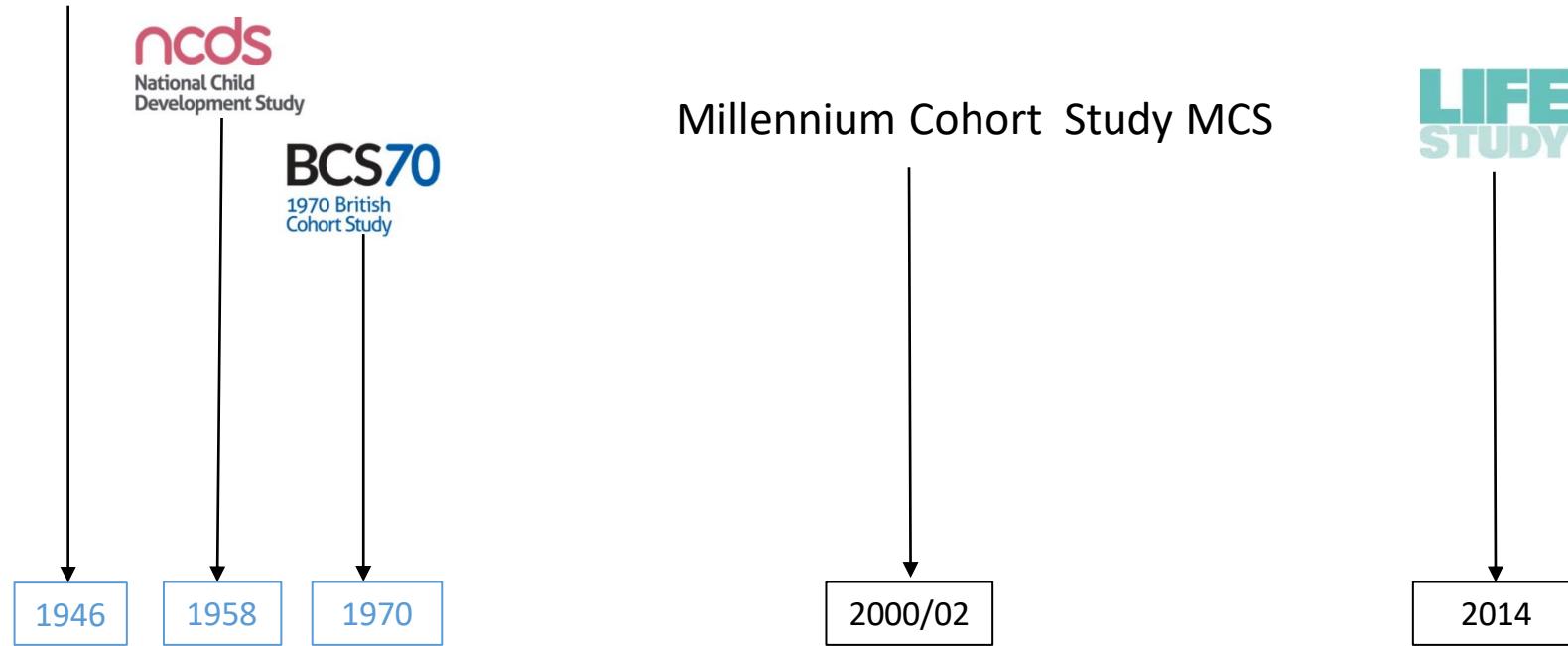
# UK Birth Cohorts

National Survey of Health  
and Development 1946



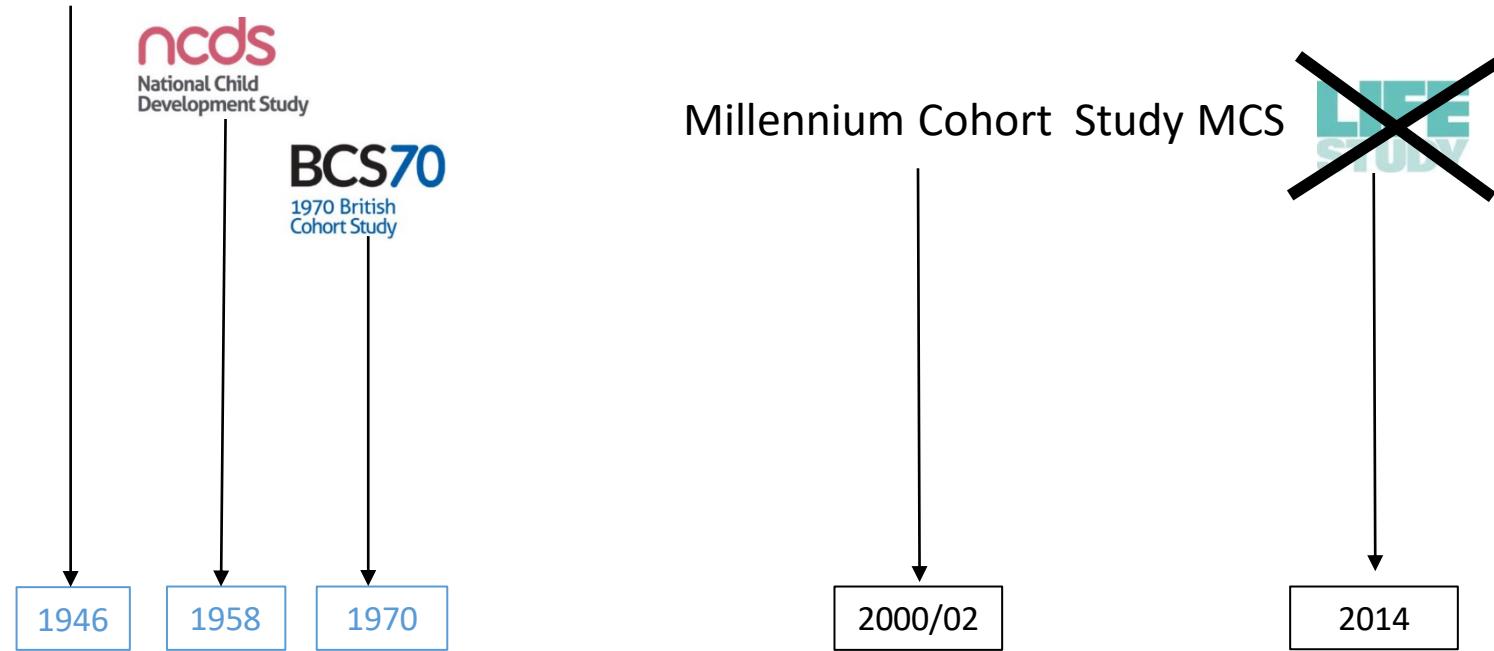
# UK Birth Cohorts

National Survey of Health  
and Development 1946



# UK Birth Cohorts

National Survey of Health  
and Development 1946





Life Study aimed to find associations between factors early in life and outcomes later on.

EPIDEMIOLOGY

# Massive UK baby study cancelled

*After demise of similar US project, decision prompts rethink about design of future cradle-to-grave efforts.*

BY HELEN PEARSON

**A**n ambitious study that planned to collect information on 80,000 British babies throughout their lives has ended just 8 months after its official launch because not enough prospective parents signed up. The closure comes less than a year after the US National Institutes of Health (NIH) cancelled a similar effort to trace 100,000 children from birth, prompting fears that researchers will now shy away from proposing similar studies.

"I am afraid that the scientific community may not dare to embark on similarly ambitious cohort studies in the near future," says Camilla Stoltenberg, who heads the Norwegian Institute of Public Health in Oslo. She is responsible for a major birth-cohort study in Norway and chaired the international scientific-advisory committee to the UK project, called Life Study.

Prized by both medical researchers and social scientists, birth-cohort studies reveal associations between factors early in life, such as poverty or a mother's diet in pregnancy, and outcomes later on, ranging from diseases to cognition and earnings. Various efforts already exist around the world, but Life Study was to be one of the biggest and most ambitious yet. It got the green light in 2011 when government funding bodies, including the Economic and Social Research Council (ESRC) and the Medical Research Council, agreed to support the study with £38.4 million (US\$58.9 million) until 2019.

In January 2015, a team led by Carol Dezateux, a paediatric epidemiologist at University College London's Institute of Child Health, opened the study's first dedicated recruitment centre, on the outskirts of London. The researchers hoped to sign up as many as 16,000 prospective mothers — of a total target of 60,000 — by July 2016. Another

MASTERFILE/ORBIS

<https://www.nature.com/articles/nature.2015.18631.pdf>

# Youth Cohort Study

[Abstract](#) | [Access](#) | [Get started](#) | [FAQ](#) | [Related](#) | [Links](#) | [Search](#)

---

## SERIES ABSTRACT

The Youth Cohort Study (YCS) began in 1985. It is a major programme of longitudinal research designed to monitor the behaviour and decisions of representative samples of young people aged 16 years onwards as they make the transition from compulsory education to further or higher education, or to the labour market. The YCS tries to identify and explain the factors which influence post-16 transitions, for example, educational attainment, training opportunities, experiences at school. To date the YCS covers 13 cohorts and over 40 surveys.

## History of Next Steps (LSYPE) in England

Next Steps, formerly known as the Longitudinal Study of Young People in England (LSYPE), is a major survey of young people born between 1 September 1989 and 31 August 1990. Originally commissioned by the then Department for Children, Schools and Families (now Department for Education, DfE), the study was designed to examine key factors affecting educational progress, attainment and transitions following the end of compulsory education.

The study began in 2004 when the young people were aged 13 to 14, and included pupils in both state and independent schools, as well as Pupil Referral Units. The first survey collected information on approximately 15,770 young people.

# New Studies on the Horizon

- Generation New Era (<https://gnestudy.info/>)
- Children of the 2020s Study  
(<https://cls.ucl.ac.uk/cls-studies/children-of-the-2020s-study/>)



# Administrative Data

- ONS – Longitudinal Study (England and Wales)
- Northern Ireland Longitudinal Study (NILS)
- Scottish Longitudinal Study
  - A panel study of 274k people based on Census records
  - <http://www.ls.ac.uk/sls/>

# Panel Dataset Examples (Household Panel Studies)

- US Panel Study of Income Dynamics (PSID)
  - began in 1968
  - <http://psidonline.isr.umich.edu/>
- Germany Socio-Economic Panel (SOEP)
  - began in 1984
  - <http://www.diw.de/en/soep>
- British Household Panel Survey BHPS
  - (1991 onwards)
  - 5k households, 10k adults,
  - <http://www.iser.essex.ac.uk/survey/bhps>



# BHPS | British Household Panel Survey

## [Home](#)

The British Household Panel Survey began in 1991 and is a multi-purpose study whose unique value resides in the fact that:

- it follows the same representative sample of individuals - the panel - over a period of years;
- it is household-based, interviewing every adult member of sampled households;
- it contains sufficient cases for meaningful analysis of certain groups such as the elderly or lone parent families.

The wave 1 panel consists of some 5,500 households and 10,300 individuals drawn from 250 areas of Great Britain. Additional samples of 1,500 households in each of Scotland and Wales were added to the main sample in 1999, and in 2001 a sample of 2,000 households was added in Northern Ireland, making the panel suitable for UK-wide research.

### [BHPS](#)

- [+ About](#)
- [Acquiring the data](#)
- [+ Documentation](#)
- [+ Scientific steering committee](#)
- [Quality profile](#)
- [+ Faqs](#)
- [Updates](#)
- [Publications](#)

**Nuisance calls claiming to be  
'British Household Survey'**  
We have recently received a

[Home](#) > [BHPs](#) > [Documentation](#) > [Volb](#)

# BHPs Documentation - Volume B - the Codebook

Menu

You may access the codebook material in several ways:

- by consulting the **Subject Category Thesaurus** in order to find suitable index term(s).
- by selecting a specific **Index Term**
- by viewing a list of BHPs **Record Types**
- by Wave:
  - [Wave One \(A\)](#)
  - [Wave Two \(B\)](#)
  - [Wave Three \(C\)](#)
  - [Wave Four \(D\)](#)
  - [Wave Five \(E\)](#)
  - [Wave Six \(F\)](#)
  - [Wave Seven \(G\)](#)
  - [Wave Eight \(H\)](#)
  - [Wave Nine \(I\)](#)



# *Understanding Society: the UK Household Longitudinal Study*



<http://www.understandingsociety.org.uk/>

- Understanding Society (US)
  - Also known as the UK Household Longitudinal Study (UKHLS)
- Began in January 2009
- Incorporates and extends the BHPS
- 40k UK households (4K Scottish households)
- 4k households in a special ethnic minorities sample
- Innovations include:
  - Linking to administrative data; spatial data; biometric data; qualitative data; child data (from age 10)

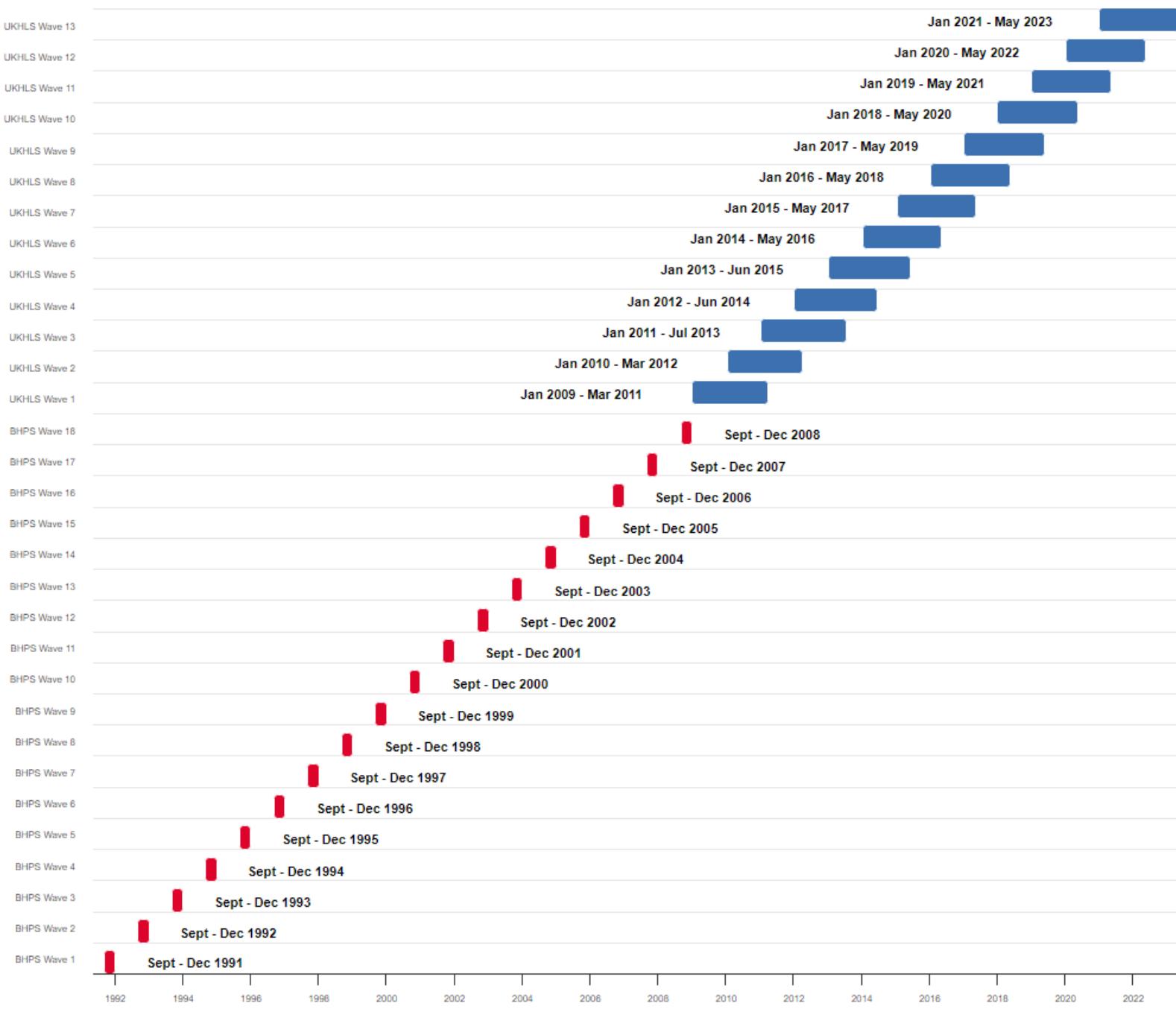
<http://www.understandingsociety.org.uk/>

# *Understanding Society Sample*

- Approx. 27,000 households -
  - The fieldwork for this sample commenced in January 2009
- A boost ethnic minority sample,
  - focussed on five main ethnic minority groups, comprising 4,000 households
- Incorporating the BHPS sample of approximately 8,400 households
- An Innovation Panel of 1500 households to enable methodological research
  - (panel began in January 2008)

# *Understanding Society*

- Focus on new research issues
- Opportunities for mixed methods:
  - Data linkage admin, organisation, spatial
  - Bio-markers and health indicators
  - Qualitative data
  - Other non-standard data: diaries, visual, audio



**Table 1: List of main data files**

Filename	Description
w_indall	Household grid data for all persons in the household, including children and non-respondents. <i>The variable pidp or the combination of variables "w_hidp w_pno" uniquely identifies each row of w_indall. The variable pidp or the combination of variables "bw_hidp bw_pno" uniquely identifies each row of bw_indall.</i>
w_hhresp	Substantive data collected from responding households. <i>The variable w_hidp uniquely identifies each row in b_hhresp. The variable bw_hidp uniquely identifies each row in bw_hhresp.</i>
bw_hhresp	
w_indresp	Substantive data collected from responding adults (16+) including proxies. Some information collected in these questionnaires are better presented in multi-level files (see Table 2). <i>The variable pidp or the combination of variables "w_hidp w_pno" uniquely identifies each row of w_indresp. The variable pidp or the combination of variables "bw_hidp bw_pno" uniquely identifies each row of bw_indresp.</i>
bw_indresp	
w_youth	Substantive data from youth questionnaire. <i>The variable pidp or the combination of variables "w_hidp w_pno" uniquely identifies each row of w_youth. The variable pidp or the combination of variables "bw_hidp bw_pno" uniquely identifies each row of bw_youth.</i>
bw_youth	
w_child	Childcare, consents and school information of all children (0-15 years) in the household. This is a derived data file collecting information pertaining to children as reported by their parents and guardians in the adult questionnaire. <i>The variable pidp or the combination of variables "w_hidp w_pno" uniquely identifies each row of w_child.</i>
bw_child	
w_egoalt	Kin and other relationships between pairs of individuals in the household. This is a derived data file based on information collected in the household grid about relationships between household members. <i>The combination of variables "pidp apidp" or "w_hidp w_pno w_apno" uniquely identifies each row in w_egoalt. The combination of variables "pidp apidp" or "bw_hidp bw_pno bw_apno" uniquely identifies each row in bw_egoalt.</i>
bw_egoalt	
w_income	This file contains reports of unearned income and state benefits for each individual. <i>The combination of variables "pidp w_fiseq" uniquely identifies each row in w_income. The combination of variables "pidp bw_ficode bw_fiseq" uniquely identifies each row in bw_income.</i>
bw_income	

## MAIN SURVEY

[Home](#) > [Data and documentation](#) > [Main survey](#) > [Variable search](#)[Variable search](#)[Datafiles](#)[Index terms](#)[Questionnaire modules](#)

# Datafiles

Find the variables you need for your research by searching by variable name or data file

Datafile	Description
<a href="#">indall</a>	Information for all persons in household, incl. children and non-respondents
<a href="#">xhhrel</a>	Family matrix
<a href="#">child</a>	Contains information about respondent's children
<a href="#">egoalt</a>	Kin and other relationships between pairs of individuals in the household
<a href="#">hhresp</a>	Substantive data from responding households
<a href="#">indresp</a>	Substantive data for responding adults (16+), incl. proxies
<a href="#">youth</a>	Substantive data from youth questionnaire (age 11-15)
<a href="#">income</a>	Income and payment information
<a href="#">empstat</a>	Employment history
<a href="#">lifemst</a>	Contains information about employment status spells
<a href="#">lifejob</a>	Contains information about jobs held in employment spells
<a href="#">jobhist</a>	Employment history
<a href="#">cohab</a>	Information about previous spells of cohabitation

Main survey Dataset docx

<https://www.understandingsociety.ac.uk/documentation/mainstage/dataset-documentation>

BIOMARKERS, GENETICS & EPIGENETICS EDUCATION EMPLOYMENT FAMILY & HOUSEHOLDS POLITICS & SOCIAL ATTITUDES More ▾

## MAIN SURVEY

Home > Documentation > Main survey > Dataset documentation

# Dataset documentation

Understanding Society collects information from everyone aged 10 and over. Around 40,000 UK households have contributed to the Study. For more background information read our [About Understanding Society Guide](#).

If you are new to using Understanding Society our [Getting Started Guide](#) will help you start working with the dataset. Our list of [key variables for the analysis of individual response data](#) may also help you.

search Mainstage Variable ▾ SEARCH

DATAFILES QUESTIONNAIRE MODULES

ALL WAVES (1991-2017)

- XWAVE (1991-2017)
- WAVE 7 (2015-2017)
- WAVE 6 (2014-2016)
- WAVE 5 (2012-2015)
- WAVE 4 (2012-2014)
- WAVE 3 (2011-2013)
- WAVE 2 (2010-2012)
- WAVE 1 (2009-2011)
- WAVE B18 (2008)
- WAVE B17 (2007)

All Waves (284 Datasets)

Dataset	Description	UKHLS	BHPS
indall	Information for all persons in household, incl. children and non-respondents	<a href="#">1</a> , <a href="#">2</a> , <a href="#">3</a> , <a href="#">4</a> , <a href="#">5</a> , <a href="#">6</a> , <a href="#">7</a>	<a href="#">B1</a> , <a href="#">B2</a> , <a href="#">B3</a> , <a href="#">B4</a> , <a href="#">B5</a> , <a href="#">B6</a> , <a href="#">B7</a> , <a href="#">B8</a> , <a href="#">B9</a> , <a href="#">B10</a> , <a href="#">B11</a> , <a href="#">B12</a> , <a href="#">B13</a> , <a href="#">B14</a> , <a href="#">B15</a> , <a href="#">B16</a> , <a href="#">B17</a> , <a href="#">B18</a>
egoalt	Kin and other relationships between pairs of individuals in the household	<a href="#">1</a> , <a href="#">2</a> , <a href="#">3</a> , <a href="#">4</a> , <a href="#">5</a> , <a href="#">6</a> , <a href="#">7</a>	<a href="#">B1</a> , <a href="#">B2</a> , <a href="#">B3</a> , <a href="#">B4</a> , <a href="#">B5</a> , <a href="#">B6</a> , <a href="#">B7</a> , <a href="#">B8</a> , <a href="#">B9</a> , <a href="#">B10</a> , <a href="#">B11</a> , <a href="#">B12</a> , <a href="#">B13</a> , <a href="#">B14</a> , <a href="#">B15</a> , <a href="#">B16</a> , <a href="#">B17</a> , <a href="#">B18</a>
hhresp	Substantive data from responding households	<a href="#">1</a> , <a href="#">2</a> , <a href="#">3</a> , <a href="#">4</a> , <a href="#">5</a> , <a href="#">6</a> , <a href="#">7</a>	<a href="#">B1</a> , <a href="#">B2</a> , <a href="#">B3</a> , <a href="#">B4</a> , <a href="#">B5</a> , <a href="#">B6</a> , <a href="#">B7</a> , <a href="#">B8</a> , <a href="#">B9</a> , <a href="#">B10</a> , <a href="#">B11</a> , <a href="#">B12</a> , <a href="#">B13</a> , <a href="#">B14</a> , <a href="#">B15</a> , <a href="#">B16</a> , <a href="#">B17</a> , <a href="#">B18</a>
indresp	Substantive data for responding adults (16+),	<a href="#">1</a> , <a href="#">2</a> , <a href="#">3</a> , <a href="#">4</a> , <a href="#">5</a> , <a href="#">6</a> , <a href="#">7</a>	<a href="#">B1</a> , <a href="#">B2</a> , <a href="#">B3</a> , <a href="#">B4</a> , <a href="#">B5</a> , <a href="#">B6</a> , <a href="#">B7</a> , <a href="#">B8</a> , <a href="#">B9</a> , <a href="#">B10</a> , <a href="#">B11</a> , <a href="#">B12</a> , <a href="#">B13</a>

OVERVIEW  
SURVEY TIMELINE  
USER GUIDES  
DATASET DOCUMENTATION  
INDEX TERMS  
QUESTIONNAIRES  
TECHNICAL REPORTS  
FIELDWORK DOCUMENTS  
LONG TERM CONTENT PLAN  
QUALITY PROFILE

Older version of the website

# Repeated Cross-Sectional Analyses

Vernon Gayle  
Professor of Sociology & Social Statistics  
University of Edinburgh

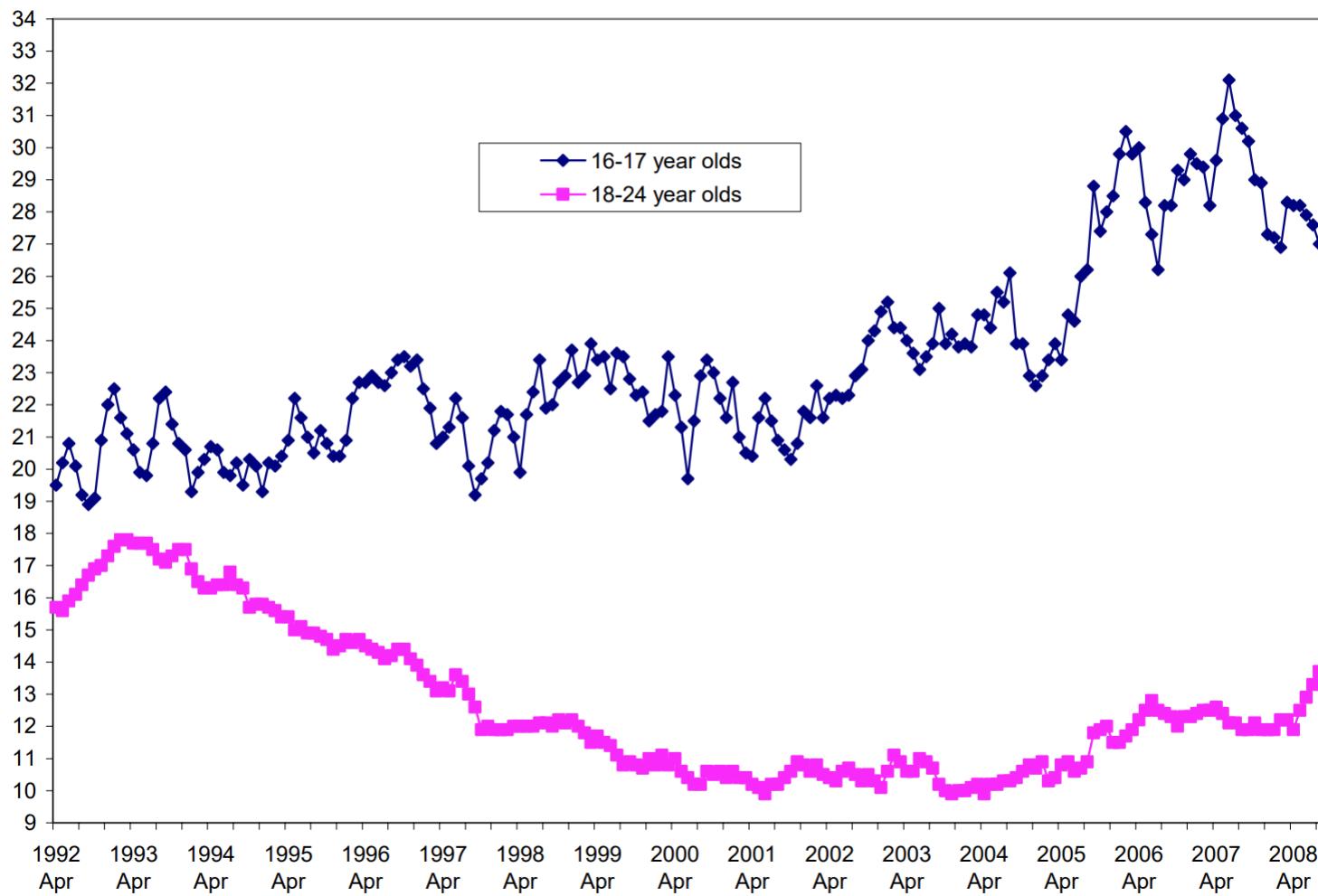
vernon.gayle@ed.ac.uk  
@profbigvern  
2024

© Vernon Gayle

# Repeated Cross-Sectional Surveys

- Often over-looked as a source of longitudinal information
- Many countries have cross-sectional surveys that are carried out on a regular basis
- They offer the possibility of pooling data for different years
- Not based on repeated contacts with the same individuals or households
- But offer opportunities to analyse general trends over time

**Figure 2. Youth unemployment, 1992-2008**



<https://www.econstor.eu/bitstream/10419/35409/1/595037216.pdf>

*Suggested Citation:* Bell, David N.F.; Blanchflower, David G. (2009) : What should be done about rising unemployment in the UK?, IZA Discussion Papers, No. 4040, Institute for the Study of Labor (IZA), Bonn,  
<https://nbn-resolving.de/urn:nbn:de:101:1-20090306334>

TABLE 2. Class and vote 1964-87

	Salariat %	Routine non-manual %	Petty bourgeoisie %	Foremen & technicians %	Working class %
<b>1964:</b>					
Con	62	58	75	38	25
Lab	19	26	14	46	68
Lib	18	16	12	15	7
	99(268)	100(197)	101(102)	99(117)	100(691)
<b>1966:</b>					
Con	61	49	67	34	24
Lab	25	41	19	61	71
Lib	15	10	15	5	5
	101(280)	100(216)	101(102)	100(111)	100(707)
<b>1970:</b>					
Con	62	51	70	39	33
Lab	29	41	19	56	61
Lib	9	9	11	5	6
	100(288)	101(187)	100(110)	100(111)	100(603)
<b>Feb 1974:</b>					
Con	54	45	68	39	24
Lab	22	30	19	40	60
Lib	24	26	13	22	16
	100(410)	100(329)	99(164)	100(106)	100(848)
<b>Oct 1974:</b>					
Con	52	44	71	35	21
Lab	23	32	13	52	64
Lib	25	24	16	13	15
	100(399)	100(307)	101(137)	99(114)	99(785)
<b>1979:</b>					
Con	61	52	77	45	32
Lab	22	32	13	43	55
Lib	17	17	10	11	13
	100(378)	100(209)	100(130)	100(150)	100(604)
<b>1983:</b>					
Con	55	53	71	44	30
Lab	13	20	12	28	49
Lib & SDP	31	27	17	28	21
	100(793)	100(547)	100(227)	100(183)	100(1127)
<b>1987:</b>					
Con	56	52	65	39	31
Lab	15	26	16	36	48
Lib & SDP	29	23	20	24	21
	100(839)	101(576)	101(245)	99(176)	100(1024)

Source: British Election Study cross-section surveys

Evans, G., Heath, A. and Payne, C., 1991. Modelling trends in the class/party relationship 1964-87. *Electoral Studies*, 10(2), pp.99-117.

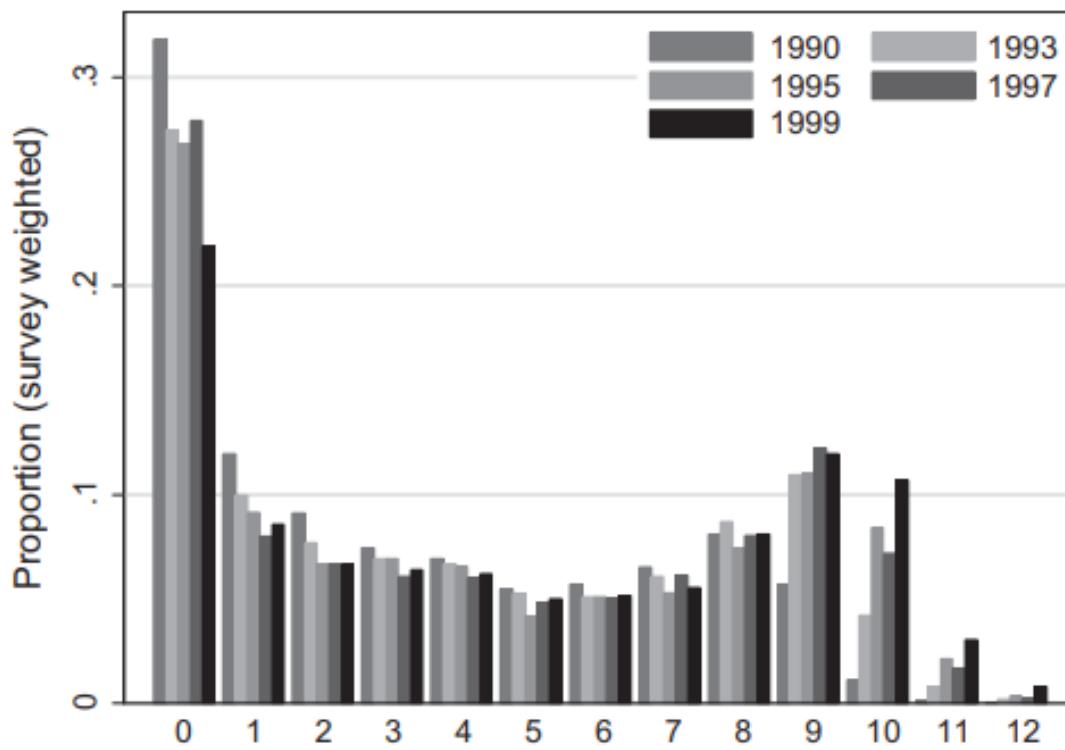


Figure 2. Number of GCSEs at grades A\*-C attained in Year 11 by school year cohort, 1990–1999.

Note: n=54,236.

# **By Slow Degrees: Two Centuries of Social Reproduction and Mobility in Britain**

**by Paul Lambert, Kenneth Prandy and Wendy Bottero**  
Stirling University

*Sociological Research Online, Volume 12, Issue 1,*  
< <http://www.socresonline.org.uk/12/1/prandy.html>>  
doi:10.5153/sro.1493

Received: 13 Nov 2006 Accepted: 19 Jan 2007 Published: 31 Jan 2007

---

## **Abstract**

This paper discusses long term trends in patterns of intergenerational social mobility in Britain. We argue that there is convincing empirical evidence of a small but steady linear trend towards increasing social mobility throughout the period 1800-2004. Our conclusions are based upon the construction and analysis of an extended micro-social dataset, which combines records from an historical genealogical study, with responses from 31 sample surveys conducted over the period 1963-2004. There has been much previous study of trends in social mobility, and little consensus on their nature. We argue that this dissension partly results from the very slow pace of change in mobility rates, which makes the time-frame of any comparison crucial, and raises important methodological questions about how long-term change in mobility is best measured. We highlight three methodological difficulties which arise when trying to draw conclusions over mobility trends - concerning the extent of controls for life course effects; the quality of data resources; and the measurement of stratification positions. After constructing a longitudinal dataset which attempts to confront these difficulties, our analyses provide robust evidence which challenges hitherto more popular, politicised claims of declining or unchanging mobility. By contrast, our findings suggest that Britain has moved, and continues to move, steadily towards increasing equality in the relationship between occupational attainment and parental background.

## Wide format dataset – single survey contact

<code>id</code>	<code>age</code>	<code>female</code>	<code>working_hours</code>
001	20	0	37
002	30	1	39
003	40	1	45

## Wide format dataset - repeated contacts

id	female	age_1971	age_1972	age_1973	work_hours_1971	work_hours_1972	work_hours_1973
001	0	20	21	22	37	40	35
002	1	30	31	32	39	40	35
003	1	40	41	42	45	45	15

## Snapshot of long format dataset

<b>id</b>	<b>year</b>	<b>age</b>	<b>hours</b>	<b>ln_wage</b>
3	68	22	40	1.49
3	69	23	40	1.70
3	70	24	40	1.45

**id:** personal identification number

**year:** year of the survey

**age:** respondent's age in years

**hours:** number of hours per week normally worked in main job

**ln\_wage:** log of weekly wages (adjusted for inflation)

# Pooled Cross-Sectional Model

- Panels are pooled together and standard statistical model used (e.g. OLS)
- A good place to start to explore
- Results provide some initial information

# Pooled Cross-Sectional Model

- Overall limitation of the pooled cross-sectional model is that it assumes that each observation (i.e. row within a long format dataset) is independent of other observations

# Pooled Cross-Sectional Model

- With panel data we know that individual respondents contribute many times to the data (usually once per wave for many waves)
- Pooling all of the data violates the standard regression modelling assumption that each observation is independent

# Pooled Cross-Sectional Model

- In practice standard errors that are too small
- Think about what this means for significance?

# Robust Standard Errors

- Robust standard errors are sometimes known as Huber/White sandwich estimates of variance (see White, 1984, Huber, 1967)

# Analysing Panel Data (Part 1)

Vernon Gayle  
Professor of Sociology & Social Statistics  
University of Edinburgh

vernon.gayle@ed.ac.uk  
@profbigvern  
2023

© Vernon Gayle

'The making of a causal inference is not a simple affair that can be reduced to a formula applied mechanically to a set of panel data on two or more variables'

(Duncan, 1972 p.36)

<b>id</b>	<b>year</b>	<b>age</b>	<b>hours</b>	<b>ln_wage</b>
3	68	22	40	1.49
3	69	23	40	1.70
3	70	24	40	1.45

## Collapsed dataset of the mean values

<b>Id</b>	<b>year <math>\bar{x}</math></b>	<b>age <math>\bar{x}</math></b>	<b>hours <math>\bar{x}</math></b>	<b>ln_wage <math>\bar{x}</math></b>
3	60	23	40	1.55

**id:** personal identification number

**year  $\bar{x}$ :** mean year of the survey

**age  $\bar{x}$ :** mean of respondent's age in years

**hours  $\bar{x}$ :** mean number of hours per week normally worked in main job

**ln\_wage  $\bar{x}$ :** mean log of weekly wages (adjusted for inflation)

## Wide format dataset – single survey contact

<code>id</code>	<code>age</code>	<code>female</code>	<code>working_hours</code>
001	20	0	37
002	30	1	39
003	40	1	45

## Wide format dataset - repeated contacts

id	female	age_1971	age_1972	age_1973	work_hours_1971	work_hours_1972	work_hours_1973
001	0	20	21	22	37	40	35
002	1	30	31	32	39	40	35
003	1	40	41	42	45	45	15

## Snapshot of long format dataset

<b>id</b>	<b>year</b>	<b>age</b>	<b>hours</b>	<b>ln_wage</b>
3	68	22	40	1.49
3	69	23	40	1.70
3	70	24	40	1.45

**id:** personal identification number

**year:** year of the survey

**age:** respondent's age in years

**hours:** number of hours per week normally worked in main job

**ln\_wage:** log of weekly wages (adjusted for inflation)

# Pooled Cross-Sectional Model

- Panels are pooled together and standard statistical model used (e.g. OLS)
- A good place to start to explore
- Results provide some initial information

# Pooled Cross-Sectional Model

- Overall limitation of the pooled cross-sectional model is that it assumes that each observation (i.e. row within a long format dataset) is independent of other observations

# Pooled Cross-Sectional Model

- With panel data we know that individual respondents contribute many times to the data (usually once per wave for many waves)
- Pooling all of the data violates the standard regression modelling assumption that each observation is independent

# Pooled Cross-Sectional Model

- In practice standard errors that are too small
- Think about what this means for significance?

# Robust Standard Errors

- Robust standard errors are sometimes known as Huber/White sandwich estimates of variance (see White, 1984, Huber, 1967)

# Between Effects Model

Estimates a standard cross-sectional model  
on the data

A regression model with  $Y$  the mean of the  
log of weekly wages (adjusted for inflation)

$X$  vars

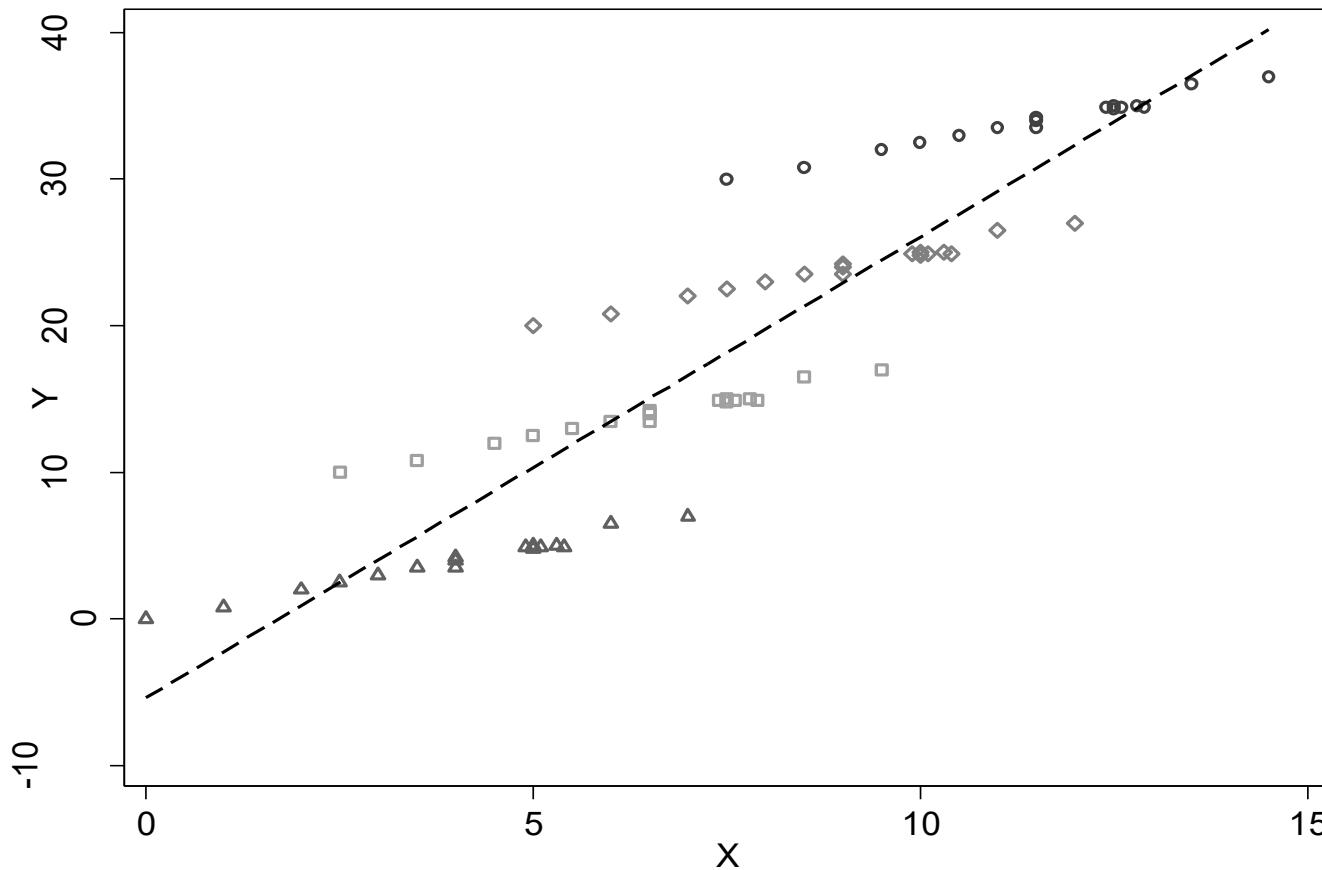
- mean hours per week normally worked in  
the respondent's main job
- mean age across the three waves of the  
survey

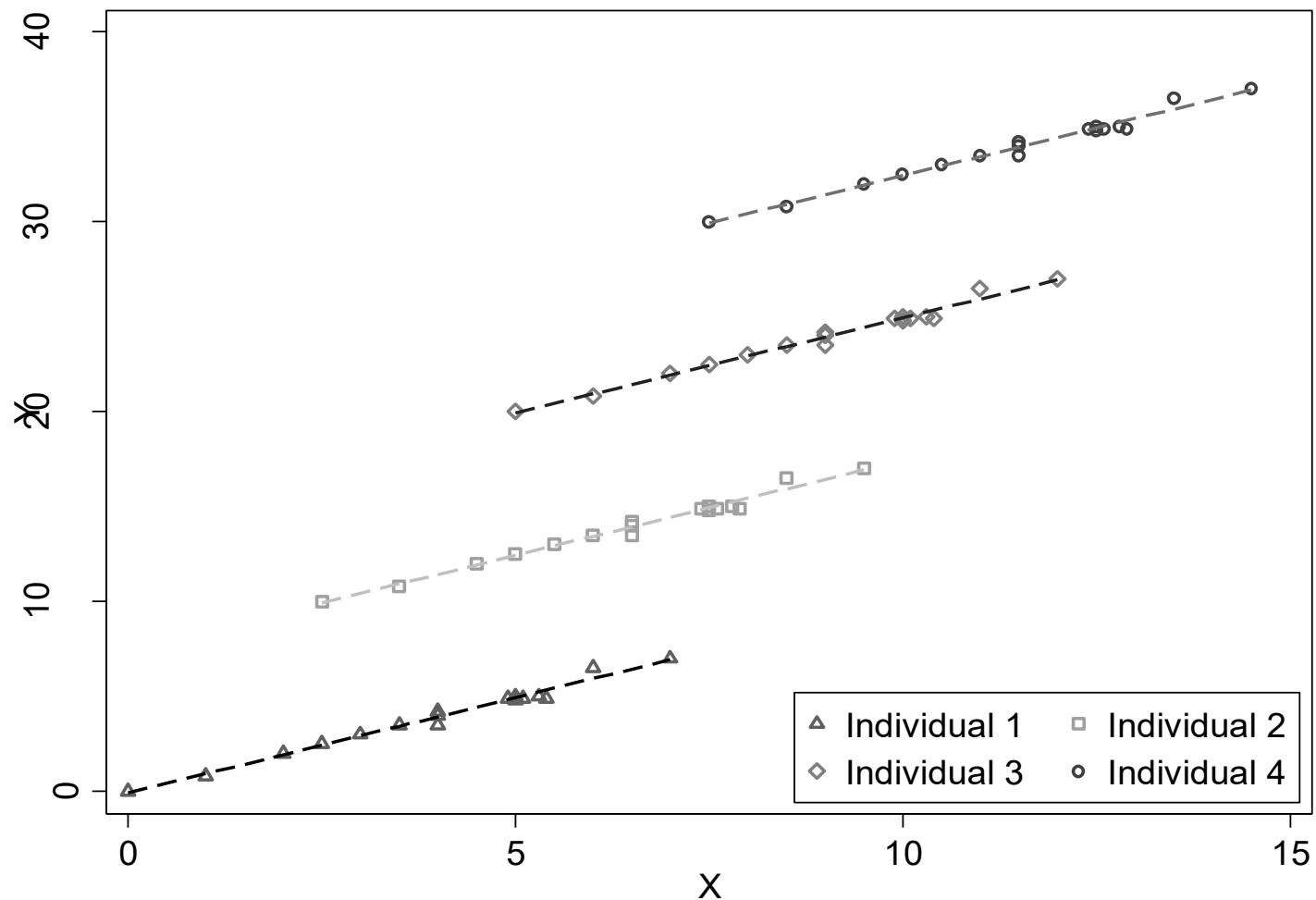
# Between Effects Model

Because now there is only one row of data per respondent the problem of non-independence of observations in the original (long format) panel data is sidestepped

*What might the limitation of this approach be?*

# A Thought Experiment...





# The Fixed Effects Panel Model

- Concentrates on change over time within an individual respondent
- Can include explanatory variables that, for the individual respondent, change over time (e.g. age, monthly income and body mass index)
- In general cannot include explanatory variables that, for the individual respondent, are time-constant (e.g. town of birth, birth weight, father's occupation when respondent was aged 14)
- Has the potentially attractive property of providing robust estimates when observed explanatory variables are correlated with the unobserved effects

# The Random Effects Panel Model

- Analyses both change within an individual respondent's outcomes, and differences between respondents' outcomes
- Can include explanatory variables that, for the individual respondent, change over time (e.g. age, monthly income and body mass index)
- Can include explanatory variables that, for the individual respondent, are time-constant (e.g. town of birth, birth weight, father's occupation when respondent was aged 14)
- Makes the assumption that observed explanatory variables are not correlated with the unobserved effects

# Notation

Pooled Cross-Sectional Regression Model

$$(1) \quad Y_{it} = \beta_0 + \beta_1 X_{1it} + \dots + \beta_k X_{kit} + \varepsilon_{it}$$

Fixed Effects Panel Regression Model

$$(2) \quad Y_{it} = \beta_0 + \lambda_i + \beta_1 X_{1it} + \dots + \beta_k X_{kit} + \varepsilon_{it}$$

Random Effects Panel Regression ('random intercepts' version)

$$(3) \quad Y_{it} = \beta_0 + \beta_1 X_{1it} + \dots + \beta_k X_{kit} + v_i + \varepsilon_{it}$$

# A Toy Example

Variable	Obs	Mean	Min	Max	Label
<hr/>					
y	32	5.38	1	11	y outcome variable
id	32	4.50	1	8	id
female	32	.50	0	1	female
wave2	32	.25	0	1	wave 2
wave3	32	.25	0	1	wave 3
wave4	32	.25	0	1	wave 4
<hr/>					

# OLS Regression

y	Coef.	Std. Err.	t	P t	[95% Conf. Interval]
-----+-----					
female	.625	.4187448	1.49	0.147	-.2341934 1.484193
wave2	.75	.5921946	1.27	0.216	-.465083 1.965083
wave3	3.5	.5921946	5.91	0.000	2.284917 4.715083
wave4	6.25	.5921946	10.55	0.000	5.034917 7.465083
_cons	2.4375	.4681709	5.21	0.000	1.476893 3.398107
-----+-----					

# Between Effects Model

Between regression (regression on group means) Number of obs = 32  
Group variable: id Number of groups = 8

R-sq:  
within = . min = 4  
between = 0.2500 avg = 4.0  
overall = 0.0133 max = 4

F(1, 6) = 2.00  
sd(u\_i + avg(e\_i.))= .625 Prob F = 0.2070

---

y	Coef.	Std. Err.	t	P t	[95% Conf. Interval]
female	.625	.4419417	1.41	0.207	-.4563925 1.706392
wave2	0	(omitted)			
wave3	0	(omitted)			
wave4	0	(omitted)			
_cons	5.0625	.3125	16.20	0.000	4.29784 5.82716

---

# Fixed Effects Model

```
Fixed-effects (within) regression                                Number of obs     =      32
Group variable: id                                         Number of groups  =       8

R-sq:                                                               Obs per group:
within  = 0.8722                                              min =         4
between = .                                                 avg =      4.0
overall = 0.8259                                              max =         4

F(3, 21) = 47.77
corr(u_i, Xb) = -0.0000
Prob F = 0.0000

-----
          y |      Coef.    Std. Err.      t    P|t|   [95% Conf. Interval]
-----+
female |          0  (omitted)
wave2 |      .75    .5824824    1.29    0.212   -.4613384    1.961338
wave3 |      3.5    .5824824    6.01    0.000    2.288662    4.711338
wave4 |      6.25    .5824824   10.73    0.000    5.038662    7.461338
_cons |      2.75    .4118772    6.68    0.000    1.893454    3.606546
-----+
sigma_u |    .6681531
sigma_e |  1.1649647
rho |  .24752475  (fraction of variance due to u_i)
-----+
F test that all u_i=0: F(7, 21) = 1.32                           Prob F = 0.2914
```

# Random Effects Model

```
Random-effects GLS regression                               Number of obs     =      32
Group variable: id                                     Number of groups  =       8

R-sq:
    within  = 0.8722                                         Obs per group:
    between = 0.2500                                         min  =        4
    overall = 0.8392                                         avg  =      4.0
                                                       max  =        4

Wald chi2(4) = 145.32
corr(u_i, X) = 0 (assumed)                           Prob chi2 = 0.0000

-----+
y | Coef. Std. Err.      z   P|z| [95% Conf. Interval]
-----+
female | .625  .4419417   1.41  0.157  -.2411899  1.49119
wave2 | .75   .5824824   1.29  0.198  -.3916445  1.891644
wave3 | 3.5   .5824824   6.01  0.000  2.358356   4.641644
wave4 | 6.25  .5824824  10.73  0.000  5.108356   7.391644
_cons | 2.4375 .474224   5.14  0.000  1.508038   3.366962
-----+
sigma_u | .22658174
sigma_e | 1.1649647
rho | .03645008 (fraction of variance due to u_i)
-----+
```

# Comparing Models

	BE	FE	RE
	b (se)	b (se)	b (se)
-----			
female	0.625 (0.442)	0.000 (..)	0.625 (0.442)
wave2	0.000 (..)	0.750 (0.582)	0.750 (0.582)
wave3	0.000 (..)	3.500*** (0.582)	3.500*** (0.582)
wave4	0.000 (..)	6.250*** (0.582)	6.250*** (0.582)
_cons	5.063*** (0.313)	2.750*** (0.412)	2.438*** (0.474)
-----			
n	32	32	32
-----			
BE:	Between Effects Model		
FE:	Fixed Effects Panel Model		
RE:	Random Effects Panel Model		

# Concluding Thoughts

Random effects panel model is using (or borrowing) some information from the fixed effects panel model

At the same time it is borrowing some information from the between effects model

This should illustrate why econometricians often make oral statements such as ‘the random effects panel model is a matrix weighted average of the within-effects (fixed effects) and the between effects’

# An Example from the British Household Panel Survey

# Example: A BHPS Panel Data File

first 10 years;  
25 - 35 year olds;  
'essex originals';  
males;  
working full-time;

yvar is wPAYNU2 Usual net pay per month: current job (now adjusted for inflation)

Variable	Obs	Unique	Mean	Min	Max	Label
pid	8412	1324	1.99e+07	1.00e+07	1.07e+08	cross-wave person identifier
wave	8412	10	5.382668	1	10	wave of the BHPS
zhid	8412	8332	5794920	1000381	1.07e+07	household identification number
zlineno	8412	7	1.497147	1	7	person number
zdobyr	8412	9	1960.981	1957	1965	year of birth
zpaynu2	6098	3898	1050.62	50.14124	9741.704	(deflated 1991)usual net pay per month...
zjbhrs	6435	60	40.47677	0	99	no. of hours normally worked per week
zjbcssm	7503	558	34.75936	.56	90.32	cambridge scale males: present job
pacssm	6827	347	30.12809	.56	85.04	cambridge scale males : father's job
graduate	7924	2	.1680969	0	1	Graduates (zqfachi)
zregage	8412	19	9.178673	0	18	age at interview-25

# Summary Statistics

Variable	Obs	Mean	Std. Dev.	Min	Max
-----+-----					
pid	8,412	19933161	15557954	10004531	107127271
wave	8,412	5.382668	2.883692	1	10
zhid	8,412	5794920	2814730	1000381	10677259
zpno	8,412	1.497147	.8551585	1	7
zdoby	8,412	1960.981	2.607186	1957	1965
-----+-----					
zpaynu2	6,098	1050.62	488.2907	50.14124	9741.704
zjbhrs	6,435	40.47677	7.895388	0	99
zjbcssm	7,503	34.75936	19.10313	.56	90.32
pacssm	6,827	30.12809	19.00986	.56	85.04
graduate	7,924	.1680969	.3739759	0	1
-----+-----					
zregage	8,412	9.178673	3.853988	0	18

# Pooled Cross-Sectional Model

```
. reg zpaynu2 zjbhr zjbcssm pacssm graduate zregage i.wav
```

Source	SS	df	MS	Number of obs	=	5,097
				F(14, 5082)	=	128.85
Model	325506963	14	23250497.3	Prob > F	=	0.0000
Residual	917001747	5,082	180441.115	R-squared	=	0.2620
				Adj R-squared	=	0.2599
Total	1.2425e+09	5,096	243820.391	Root MSE	=	424.78

# Pooled Cross-Sectional Model (Continued)

		Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
zpaynu2						
zjbhrs		9.559008	.8714275	10.97	0.000	7.850635 11.26738
zjbcssm		8.709623	.370681	23.50	0.000	7.982928 9.436317
pacssm		2.099019	.354503	5.92	0.000	1.40404 2.793998
graduate		195.6784	17.63807	11.09	0.000	161.1001 230.2566
zregage		6.51771	2.257654	2.89	0.004	2.091734 10.94368
wave						
2		32.02204	26.09955	1.23	0.220	-19.14433 83.1884
3		47.53042	26.53109	1.79	0.073	-4.481953 99.5428
4		32.50858	27.03488	1.20	0.229	-20.49144 85.5086
5		37.20393	27.6692	1.34	0.179	-17.03963 91.44749
6		83.98088	28.31416	2.97	0.003	28.47293 139.4888
7		72.86194	28.92154	2.52	0.012	16.16326 129.5606
8		96.8642	29.82545	3.25	0.001	38.39346 155.3349
9		139.523	31.31379	4.46	0.000	78.13451 200.9116
10		138.2166	32.437	4.26	0.000	74.62611 201.8071
_cons		135.0271	43.56656	3.10	0.002	49.61787 220.4363

# Pooled Cross-Sectional Model With robust standard errors

```
. reg zpaynu2 zjbhhr zjbcssm pacssm graduate zregage i.wav, cluster(pid)
```

Linear regression	Number of obs	=	5,097
	F(14, 824)	=	23.09
	Prob > F	=	0.0000
	R-squared	=	0.2620
	Root MSE	=	424.78

(Std. Err. adjusted for 825 clusters in pid)

---

# Pooled Cross-Sectional Model With robust standard errors

Robust						
zpaynu2		Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
zjbhrs		9.559008	1.943735	4.92	0.000	5.743753 13.37426
zjbcssm		8.709623	.7934006	10.98	0.000	7.152299 10.26695
pacssm		2.099019	1.115752	1.88	0.060	-.0910308 4.289069
graduate		195.6784	59.72225	3.28	0.001	78.45271 312.904
zregage		6.51771	4.515835	1.44	0.149	-2.346185 15.3816
wave						
2		32.02204	13.90313	2.30	0.022	4.732323 59.31175
3		47.53042	18.75842	2.53	0.011	10.71052 84.35033
4		32.50858	20.47166	1.59	0.113	-7.674154 72.69131
5		37.20393	25.89333	1.44	0.151	-13.62072 88.02857
6		83.98088	30.03699	2.80	0.005	25.02286 142.9389
7		72.86194	34.50542	2.11	0.035	5.133077 140.5908
8		96.8642	37.73043	2.57	0.010	22.80514 170.9233
9		139.523	43.37522	3.22	0.001	54.3841 224.662
10		138.2166	46.56261	2.97	0.003	46.82133 229.6119
_cons		135.0271	76.33256	1.77	0.077	-14.80204 284.8562

# The Fixed Effects Panel Model

- Concentrates on change over time within an individual respondent
- Can include explanatory variables that, for the individual respondent, change over time (e.g. age, monthly income and body mass index)
- In general cannot include explanatory variables that, for the individual respondent, are time-constant (e.g. town of birth, birth weight, father's occupation when respondent was aged 14)
- Has the potentially attractive property of providing robust estimates when observed explanatory variables are correlated with the unobserved effects

# The Random Effects Panel Model

- Analyses both change within an individual respondent's outcomes, and differences between respondents' outcomes
- Can include explanatory variables that, for the individual respondent, change over time (e.g. age, monthly income and body mass index)
- Can include explanatory variables that, for the individual respondent, are time-constant (e.g. town of birth, birth weight, father's occupation when respondent was aged 14)
- Makes the assumption that observed explanatory variables are not correlated with the unobserved effects

# Notation

Pooled Cross-Sectional Regression Model

$$(1) \quad Y_{it} = \beta_0 + \beta_1 X_{1it} + \dots + \beta_k X_{kit} + \varepsilon_{it}$$

Fixed Effects Panel Regression Model

$$(2) \quad Y_{it} = \beta_0 + \lambda_i + \beta_1 X_{1it} + \dots + \beta_k X_{kit} + \varepsilon_{it}$$

Random Effects Panel Regression ('random intercepts' version)

$$(3) \quad Y_{it} = \beta_0 + \beta_1 X_{1it} + \dots + \beta_k X_{kit} + v_i + \varepsilon_{it}$$

# Analysing Panel Data (Part 2)

Vernon Gayle  
Professor of Sociology & Social Statistics  
University of Edinburgh

vernon.gayle@ed.ac.uk  
@profbigvern  
2024

© Vernon Gayle

# Notation

Pooled Cross-Sectional Regression Model

$$(1) \quad Y_{it} = \beta_0 + \beta_1 X_{1it} + \dots + \beta_k X_{kit} + \varepsilon_{it}$$

Fixed Effects Panel Regression Model

$$(2) \quad Y_{it} = \beta_0 + \lambda_i + \beta_1 X_{1it} + \dots + \beta_k X_{kit} + \varepsilon_{it}$$

Random Effects Panel Regression ('random intercepts' version)

$$(3) \quad Y_{it} = \beta_0 + \beta_1 X_{1it} + \dots + \beta_k X_{kit} + v_i + \varepsilon_{it}$$

```
. xtreg zpaynu2 zjbhr zjbcssm pacssm graduate zregage i.wav, fe
```

note: pacssm omitted because of collinearity

Fixed-effects (within) regression

Number of obs = 5,097

Group variable: pid

Number of groups = 825

R-sq:

within = 0.0973

Obs per group:

min = 1

between = 0.0047

avg = 6.2

overall = 0.0151

max = 10

F(13, 4259) = 35.30

corr(u\_i, Xb) = -0.0716

Prob > F = 0.0000

```
. xtreg zpaynu2 zjbhr zjbcssm pacssm graduate zregage i.wav, fe
```

```
note: pacssm omitted because of collinearity
```

Fixed-effects (within) regression

Number of obs = 5,097

Group variable: pid

Number of groups = 825

R-sq:

Obs per group:

within = 0.0973 min = 1

between = 0.0047 avg = 6.2

overall = 0.0151 max = 10

F(13, 4259) = 35.30

corr(u\_i, Xb) = -0.0716 Prob > F = 0.0000

```
. xtreg zpaynu2 zjbhr zjbcssm pacssm graduate zregage i.wav, fe
```

note: pacssm omitted because of collinearity

Fixed-effects (within) regression

Number of obs = 5,097

Group variable: pid

Number of groups = 825

R-sq:

within = 0.0973

between = 0.0047

overall = 0.0151

Obs per group:

min = 1

avg = 6.2

max = 10

F(13, 4259) = 35.30

corr(u\_i, Xb) = -0.0716

Prob > F = 0.0000

```
. xtreg zpaynu2 zjbhr zjbcssm pacssm graduate zregage i.wav, fe
```

```
note: pacssm omitted because of collinearity
```

```
Fixed-effects (within) regression
```

```
Number of obs = 5,097
```

```
Group variable: pid
```

```
Number of groups = 825
```

R-sq:

within = 0.0973

Obs per group:

min = 1

between = 0.0047

avg = 6.2

overall = 0.0151

max = 10

corr(u\_i, Xb) = -0.0716

F(13, 4259) = 35.30

Prob > F = 0.0000

zpaynu2		Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
zjbhrs		3.162216	.7598803	4.16	0.000	1.672455 4.651978
zjbcssm		.7083871	.4688331	1.51	0.131	-.2107702 1.627544
pacssm		0	(omitted)			
graduate		-68.42179	58.82074	-1.16	0.245	-183.7411 46.89752
zregage		9.779486	18.59062	0.53	0.599	-26.66781 46.22679
wave						
2		33.77626	23.69473	1.43	0.154	-12.67775 80.23027
3		50.54727	39.7453	1.27	0.204	-27.37422 128.4688
4		36.4279	57.69603	0.63	0.528	-76.68639 149.5422
5		38.07226	75.6673	0.50	0.615	-110.2751 186.4196
6		85.96321	93.47347	0.92	0.358	-97.2935 269.2199
7		87.99939	111.601	0.79	0.430	-130.7967 306.7955
8		115.7779	130.037	0.89	0.373	-139.1623 370.7181
9		161.5511	148.3231	1.09	0.276	-129.2394 452.3416
10		172.5317	167.0504	1.03	0.302	-154.9742 500.0375
_cons		751.7321	98.72638	7.61	0.000	558.177 945.2873

---

```
+-----  
sigma_u | 436.52027  
sigma_e | 251.17185  
rho | .75126958 (fraction of variance due to u_i)  
-----
```

---

```
F test that all u_i=0: F(824, 4259) = 12.59                      Prob > F = 0.0000
```

# areg

```
. areg zpaynu2 zjbhr zjbcssm pacssm graduate zregage i.wav, absorb(pid)
note: pacssm omitted because of collinearity
```

```
Linear regression, absorbing indicators                               Number of obs      =      5,097
                                                               F( 13,    4259)    =     35.30
                                                               Prob > F        =     0.0000
                                                               R-squared       =     0.7838
                                                               Adj R-squared   =     0.7413
                                                               Root MSE        =   251.1718
```

		Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
zpaynu2						
zjbhrs		3.162216	.7598803	4.16	0.000	1.672455 4.651978
zjbcssm		.7083871	.4688331	1.51	0.131	-.2107702 1.627544
pacssm		0	(omitted)			
graduate		-68.42179	58.82074	-1.16	0.245	-183.7411 46.89752
zregage		9.779486	18.59062	0.53	0.599	-26.66781 46.22679
wave						
2		33.77626	23.69473	1.43	0.154	-12.67775 80.23027
3		50.54727	39.7453	1.27	0.204	-27.37422 128.4688
4		36.4279	57.69603	0.63	0.528	-76.68639 149.5422
5		38.07226	75.6673	0.50	0.615	-110.2751 186.4196
6		85.96321	93.47347	0.92	0.358	-97.2935 269.2199
7		87.99939	111.601	0.79	0.430	-130.7967 306.7955
8		115.7779	130.037	0.89	0.373	-139.1623 370.7181
9		161.5511	148.3231	1.09	0.276	-129.2394 452.3416
10		172.5317	167.0504	1.03	0.302	-154.9742 500.0375
_cons		751.7321	98.72638	7.61	0.000	558.177 945.2873
pid		F(824, 4259) =	12.593	0.000		(825 categories)

# xtreg , re (random effects)

```
. xtreg zpaynu2 zjbhr zjbcssm pacssm graduate zregage i.wav, re
```

Random-effects GLS regression  
Group variable: pid

R-sq:

within = 0.0875  
between = 0.2458  
overall = 0.2291

corr(u\_i, X) = 0 (assumed)

Number of obs = 5,097  
Number of groups = 825

Obs per group:

min = 1  
avg = 6.2  
max = 10

Wald chi2(14) = 688.95  
Prob > chi2 = 0.0000

# xtreg , re (random effects)

```
. xtreg zpaynu2 zjbhr zjbcssm pacssm graduate zregage i.wav, re
```

Random-effects GLS regression  
Group variable: pid

R-sq:

within = 0.0875  
between = 0.2458  
overall = 0.2291

corr(u\_i, X) = 0 (assumed)

Number of obs = 5,097  
Number of groups = 825

Obs per group:

min = 1  
avg = 6.2  
max = 10

Wald chi2(14) = 688.95  
Prob > chi2 = 0.0000

# xtreg , re (random effects)

```
. xtreg zpaynu2 zjbhr zjbcssm pacssm graduate zregage i.wav, re
```

Random-effects GLS regression  
Group variable: pid

Number of obs = 5,097  
Number of groups = 825

R-sq:

within = 0.0875  
between = 0.2458  
overall = 0.2291

Obs per group:

min = 1  
avg = 6.2  
max = 10

corr(u\_i, X) = 0 (assumed)

Wald chi2(14) = 688.95  
Prob > chi2 = 0.0000

# xtreg , re (random effects)

```
. xtreg zpaynu2 zjbhr zjbcssm pacssm graduate zregage i.wav, re
```

Random-effects GLS regression  
Group variable: pid

R-sq:

within = 0.0875  
between = 0.2458  
overall = 0.2291

corr(u\_i, X) = 0 (assumed)

Number of obs = 5,097  
Number of groups = 825

Obs per group:

min = 1  
avg = 6.2  
max = 10

Wald chi2(14) = 688.95  
Prob > chi2 = 0.0000

# xtreg , re (random effects)

```
. xtreg zpaynu2 zjbhr zjbcssm pacssm graduate zregage i.wav, re
```

Random-effects GLS regression  
Group variable: pid

R-sq:  
within = 0.0875  
between = 0.2458  
overall = 0.2291

corr(u\_i, X) = 0 (assumed)

Number of obs = 5,097  
Number of groups = 825

Obs per group:  
min = 1  
avg = 6.2  
max = 10

Wald chi2(14) = 688.95  
Prob > chi2 = 0.0000

zpaynu2		Coef.	Std. Err.	z	P>  z	[ 95% Conf. Interval]
zjbhrs		4.095316	.7257752	5.64	0.000	2.672823 5.51781
zjbcssm		3.190851	.4147334	7.69	0.000	2.377988 4.003713
pacssm		3.949381	.7204518	5.48	0.000	2.537321 5.36144
graduate		202.0455	30.88336	6.54	0.000	141.5152 262.5758
zregage		9.659352	4.59312	2.10	0.035	.6570023 18.6617
wave						
2		32.94832	16.59421	1.99	0.047	.4242625 65.47238
3		48.03321	18.53548	2.59	0.010	11.70434 84.36207
4		31.86777	21.25793	1.50	0.134	-9.797011 73.53255
5		31.99669	24.47987	1.31	0.191	-15.98297 79.97635
6		78.77641	27.94861	2.82	0.005	23.99814 133.5547
7		77.89485	31.66425	2.46	0.014	15.83407 139.9556
8		100.9134	35.62999	2.83	0.005	31.07988 170.7469
9		144.7248	39.838	3.63	0.000	66.64376 222.8058
10		153.3713	44.0223	3.48	0.000	67.08919 239.6534
_cons		441.8567	45.78108	9.65	0.000	352.1274 531.5859

Random-effects GLS regression  
Group variable: pid

Number of obs = 5,097  
Number of groups = 825

R-sq:

within = 0.0875  
between = 0.2458  
overall = 0.2291

Obs per group:

min = 1  
avg = 6.2  
max = 10

corr(u\_i, X) = 0 (assumed)

Wald chi2(14) = 688.95  
Prob > chi2 = 0.0000

	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
zpaynu2						
zjbhrs	4.095316	.7257752	5.64	0.000	2.672823	5.51781
zjbcssm	3.190851	.4147334	7.69	0.000	2.377988	4.003713
pacssm	3.949381	.7204518	5.48	0.000	2.537321	5.36144
graduate	202.0455	30.88336	6.54	0.000	141.5152	262.5758
zregage	9.659352	4.59312	2.10	0.035	.6570023	18.6617
wave						
2	32.94832	16.59421	1.99	0.047	.4242625	65.47238
3	48.03321	18.53548	2.59	0.010	11.70434	84.36207
4	31.86777	21.25793	1.50	0.134	-9.797011	73.53255
5	31.99669	24.47987	1.31	0.191	-15.98297	79.97635
6	78.77641	27.94861	2.82	0.005	23.99814	133.5547
7	77.89485	31.66425	2.46	0.014	15.83407	139.9556
8	100.9134	35.62999	2.83	0.005	31.07988	170.7469
9	144.7248	39.838	3.63	0.000	66.64376	222.8058
10	153.3713	44.0223	3.48	0.000	67.08919	239.6534
_cons	441.8567	45.78108	9.65	0.000	352.1274	531.5859
sigma_u	331.89925					
sigma_e	251.17185					
rho	.63584801	(fraction of variance due to u_i)				

The total error variance can be considered as

$$(\sigma_u)^2 + (\sigma_e)^2$$

# Rho (random effects)

the fraction of the variance that is due to  $u_i$  is

$$(\sigma_u)^2$$

---

$$\frac{(\sigma_u)^2}{((\sigma_u)^2 + (\sigma_e)^2)}$$

Rho in this example...

$$(331.89925^2)/((331.89925^2)+(251.17185^2))$$

Rho = 0.64

64 per cent of the error variance is at the panel level

# Rho in this example...

$$(331.89925^2)/((331.89925^2)+(251.17185^2))$$

Rho = 0.64

64 per cent of the error variance is at the panel level

*Rho is analogous to an intra-class correlation, or ICC, as it is known in other areas such as the literature on multilevel modelling*

# Rho

When rho is zero the panel-level variance component is unimportant and the panel estimator is not different from the pooled (i.e. cross-sectional) estimator

In our experience it is almost never the case that rho is zero when estimating a model using a genuine panel dataset

One notable exception is reported in Exeter (2004), where the units in the panel were geographical areas

# Formal test of random effects

Breusch and Pagan Lagrangian multiplier test for random effects

zpaynu2[pid,t] = Xb + u[pid] + e[pid,t]

Estimated results:

	Var	sd = sqrt(Var)
zpaynu2	243820.4	493.7817
e	63087.3	251.1718
u	110157.1	331.8993

Test: Var(u) = 0  
chibar2(01) = 6525.74  
Prob > chibar2 = 0.0000

# Which Model Should We Use?

Gelman and Hill (2007), two leading statisticians, comment that the statistical literature is full of confusing and contradictory advice

Searle, Casella and McCulloch (1992) assert that because of conflicting definitions, it is no surprise that clear answers to the question ‘fixed or random effects?’ are unusual

- Researchers routinely ask, ‘Should I choose a fixed effects panel model or a random effects panel model?’
- The answer depends on what the data analyst is attempting to model
- The fixed effects panel model focuses upon the within-subject change
- The random effects panel model is influenced by both within-subject and between-subject patterns
- Both have potential advantages and limitations

```
.hausman fe re
```

----- Coefficients -----

	(b)	(B)	(b-B)	sqrt(diag(V_b-V_B))
	fe	re	Difference	S.E.
-----				
logq	.9192846	.9066805	.0126041	.0153877
logf	.4174918	.4227784	-.0052867	.0058583
lf	-1.070396	-1.064499	-.0058974	.0255088

-----

b = consistent under Ho and Ha; obtained from xtreg

B = inconsistent under Ha, efficient under Ho; obtained from xtreg

Test: Ho: difference in coefficients not systematic

$$\begin{aligned} \text{chi2}(3) &= (\mathbf{b}-\mathbf{B})' [(\mathbf{V}_b-\mathbf{V}_B)^{-1}] (\mathbf{b}-\mathbf{B}) \\ &= 2.12 \end{aligned}$$

Prob>chi2 = 0.5469

( $\mathbf{V}_b-\mathbf{V}_B$  is not positive definite)

# Which Model Should We Use?

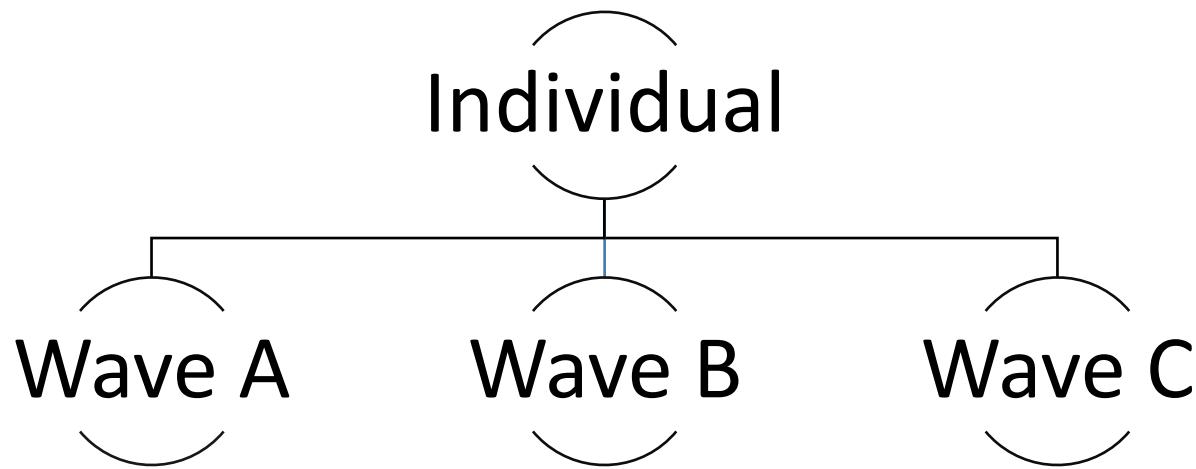
See Gayle and Lambert (2021)

*Comparing fixed effects panel models and random effects panel models (page 86-95)*

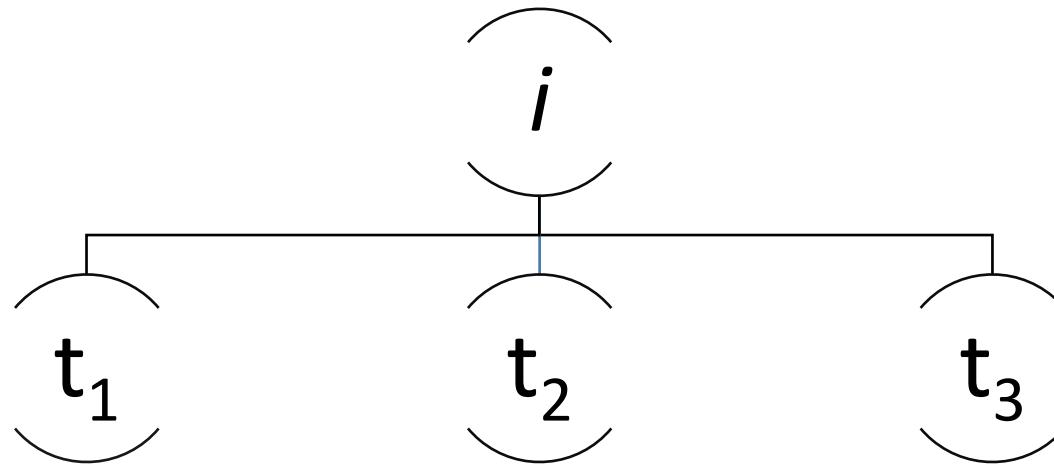
# Alternative Literatures

- Econometrics
- Multilevel modelling (e.g. in the education literature)
- Epidemiology and biostatistics (e.g. public health)

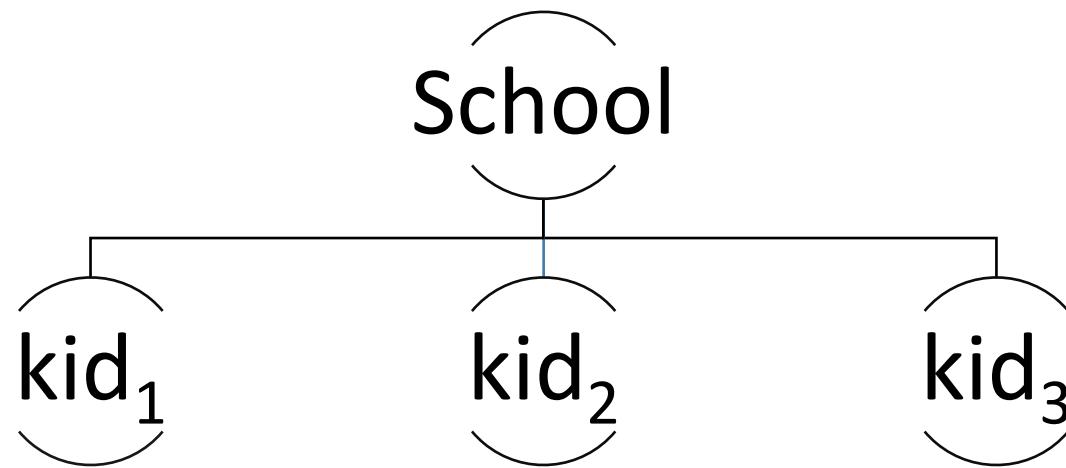
# Panel Data Structure



# Panel Data Structure



# Hierarchical Data Structure



	(1)	(2)
	Panel RE	Multilevel
Zpaynu2		
zjbhrs	3.983*** (0.723)	3.983*** (0.721)
zjbcssm	2.975*** (0.425)	2.975*** (0.415)
pacssm	4.075*** (0.759)	4.075*** (0.757)
graduate	195.5*** (32.05)	195.5*** (31.91)
zregage	9.742* (4.820)	9.742* (4.820)

_cons	449.7***	449.7***
	(46.94)	(46.84)
<hr/>		
sigma_u		
_cons	358.2***	
	(10.18)	
<hr/>		
sigma_e		
_cons	251.9***	
	(2.734)	
<hr/>		
lns1_1_1		
_cons		5.881***
		(0.0284)
<hr/>		
lnsig_e		
_cons		5.529***
		(0.0109)
<hr/>		
N	5097	5097
<hr/>		

Standard errors in parentheses

\* p<0.05, \*\* p<0.01, \*\*\* p<0.001

Stata logs the sigma\_u and sigma\_e  
in the standard xtmixed output

# Mundlak-Chamberlain

In a nutshell...

The inclusion of means of the time-varying covariates  
within the random effects model

**Table 6** Coefficients ( $\beta$ ) and their standard errors (se) from the fixed effects panel model, the random effects panel model with Mundlak adjustment, and the random effects panel model – log wages

	(1) Fixed effects $\beta_{fe}$ (s.e.)	(2) Random effects with Mundlak $\beta_{mre}$ (s.e.)		(3) Random effects $\beta_{re}$ (s.e.)	
Ft work	0.097 (0.001)	*** (0.001)	0.097 (0.001)	*** (0.001)	0.057 (0.001)
Experience(years)					
Weeks worked	0.001 (0.001)	*	0.001 (0.001)	*	0.002 (0.001)
Blue-collar occupation	-0.021 (0.014)		-0.021 (0.014)		-0.108 (0.016)
Individual's mean ft workexperience(years)			-0.090 (0.002)		***
Individual's mean weeks worked			0.010 (0.004)		**
Individual's mean blue-collar occupation			-0.316 (0.034)		***
Constant	4.709 (0.038)	*** (0.212)	6.164 (0.212)	*** (0.047)	5.523 (0.047)
n	4165		4165		4165

# Other Panel Models

## Binary Outcomes

`xtlogit`

`xtprobit`

`clogit`

## Ordinal Outcomes

`xtologit` random-effects ordered logistic models

`xtoprobit` random-effects ordered probit models

## Count Data

`xtpoisson` panel data poisson models

`xtnbreg` panel data negative binomial models

# Dynamic Models

- Dynamic panel models extend panel models
- Appeal to the idea of using panel data to better understand ‘state dependence’
- Lagged dependent variables as X vars
- Complicated because the lagged dependent variables will themselves be influenced by unobserved effects

# Dynamic Models

- Standard panel estimation procedures will be inconsistent with lagged dependent variables
- Arellano and Bond (1991) derived a suitable estimator which is available using the Stata command *xtabond*
- Stewart (2006) *redprob*

# Further Topics to Consider...

- The estimation and interpretation of interaction effects in statistical models (see Ai and Norton, 2003; Norton, Wang and Ai, 2004; Mitchell and Chen, 2005)
- Post-estimation measures and model evaluation (see Long and Freese, 2014)
- Missing data (see Carpenter and Kenward, 2012)
- Sample attrition, panel conditioning, interviewer effects and data collection modes (see Lynn, 2009)

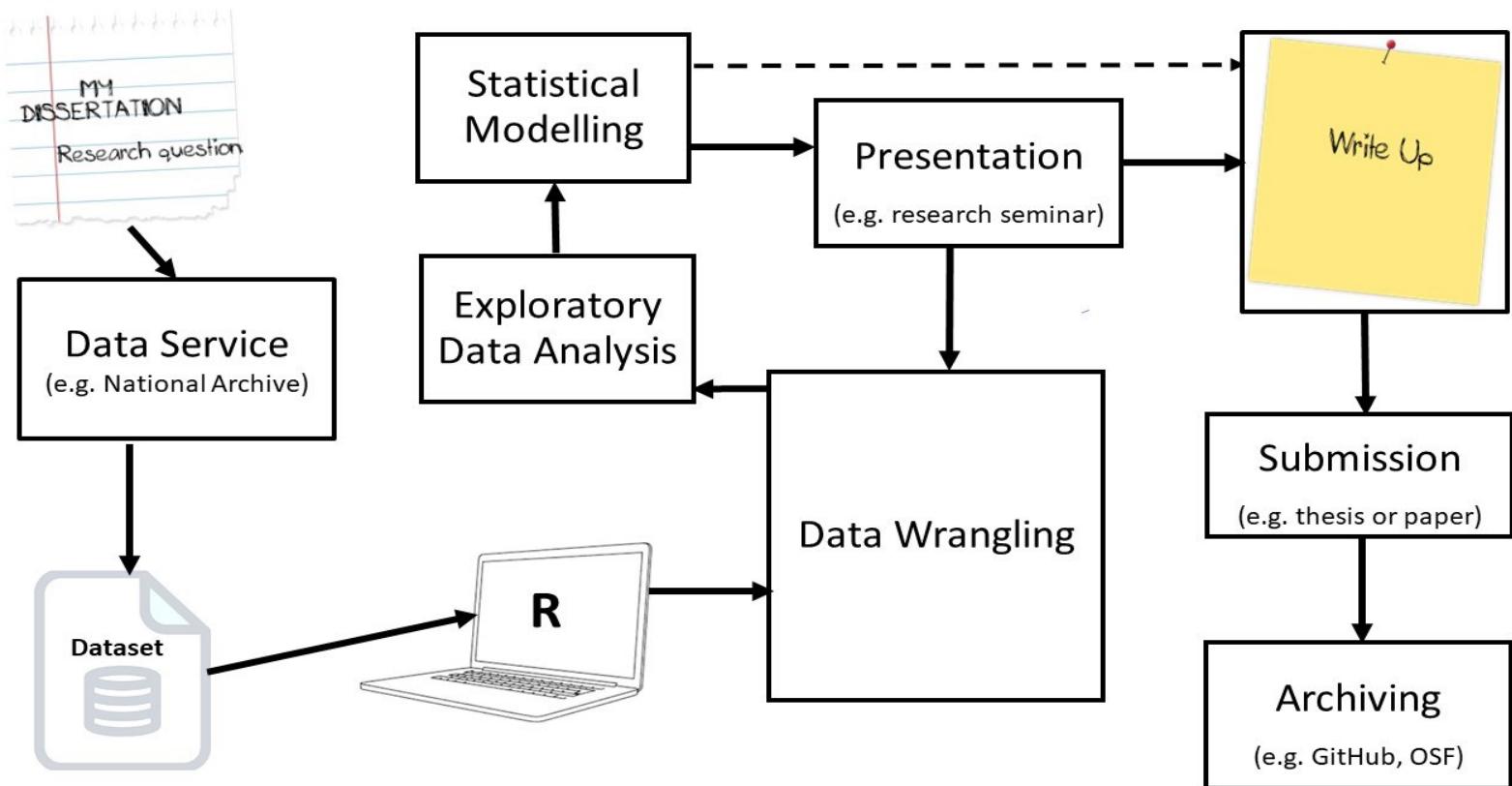
# Software and Computing

Vernon Gayle  
Professor of Sociology & Social Statistics  
University of Edinburgh

vernon.gayle@ed.ac.uk  
@profbigvern  
2024

© Vernon Gayle

# The Data Analysis Workflow



# Tools of the Trade



Stata <https://www.stata.com/>

SPSS <https://www.ibm.com/uk-en/analytics/spss-statistics-software>

SAS <https://www.sas.com/>

R <https://www.r-project.org/>

Python <https://www.python.org/>  
[\(https://www.statsmodels.org/stable/index.html\)](https://www.statsmodels.org/stable/index.html)

## Stata

```
logit admit gre gpa
```

## SPSS

```
logistic regression admit with gre gpa.
```

## SAS

```
proc logistic data="c:\data\binary" descending;  
  class rank / param=ref ;  
  model admit = gre gpa;  
run;
```

## R

```
mylogit <- glm(admit ~ gre + gpa, data = mydata, family = "binomial")
```

## Python

```
independentVar = ['gre', 'gpa', 'Int']  
logReg = sm.Logit(df['admit'], df[independentVar])  
answer = logReg.fit()
```

# Stata Code

## Panel Model

```
xtreg zpaynu2 zjbhr zjbcssm pacssm graduate zregage i.wav, re mle
```

## Multilevel Model

```
xtmixed zpaynu2 zjbhr zjbcssm pacssm graduate zregage i.wav|| pid:, mle
```

# Random Effect Panel Model R

```
library(foreign)

library(plm)

panel <- read.dta("C:\bhps_panel.dta")

remodel <- plm(zpaynu2 ~ zjbhr zjbcssm pacssm graduate
zregage, data=panel, index=c("pid", "wave"),
model="random")
```

# Considerations

1. Supervisor's expertise
2. Peer group (e.g. other PhD students)
3. Departmental access and support
4. University licenses
5. Data format and meta data (e.g. UK Data Service)
6. Academic subject area
7. Academic job market
8. Non-academic job market

# **Some Concluding Remarks**

- Some research questions require longitudinal data
- Longitudinal data are not a panacea

- For many social research projects cross-sectional data will be sufficient
- Most social research projects can be improved by the analysis of longitudinal data
- *Researchers are likely to make more rapid progress using existing large-scale longitudinal data resources*

# Final Comment...

*Angrist and Pischke (2008) playfully remarked that if applied research was easy then theorists would do it!*

*They also reassure readers that applied research is not as hard as the dense pages of Econometrica might lead us to believe*

# Longitudinal Data and Research – A Deep Dive

Vernon Gayle  
Professor of Sociology & Social Statistics  
University of Edinburgh

vernon.gayle@ed.ac.uk  
@profbigvern  
2024