# The Data Analysis Workflow
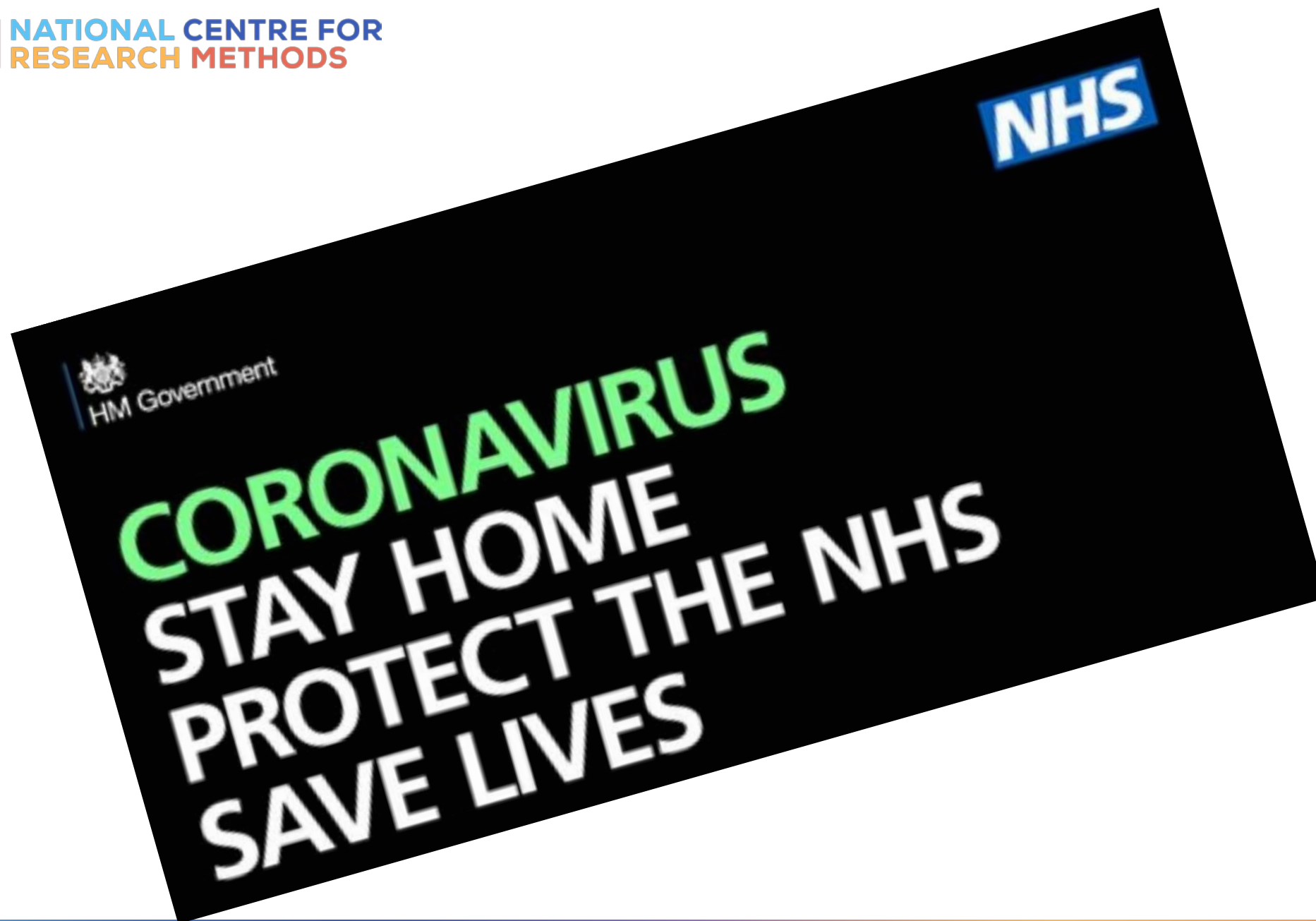
Professor Vernon Gayle
vernon.gayle@ed.ac.uk
@Profbigvern
https://github.com/vernongayle

# The Data Analysis Workflow

# A Thought Experiment (be honest…)

1. Have you ever lost a file?

2. Have you ever wondered if you have deleted a file?

3. Have you and a colleague ever been working on different versions of a file?

# A Thought Experiment (be honest…)

4. Have you ever struggled to identify data files?
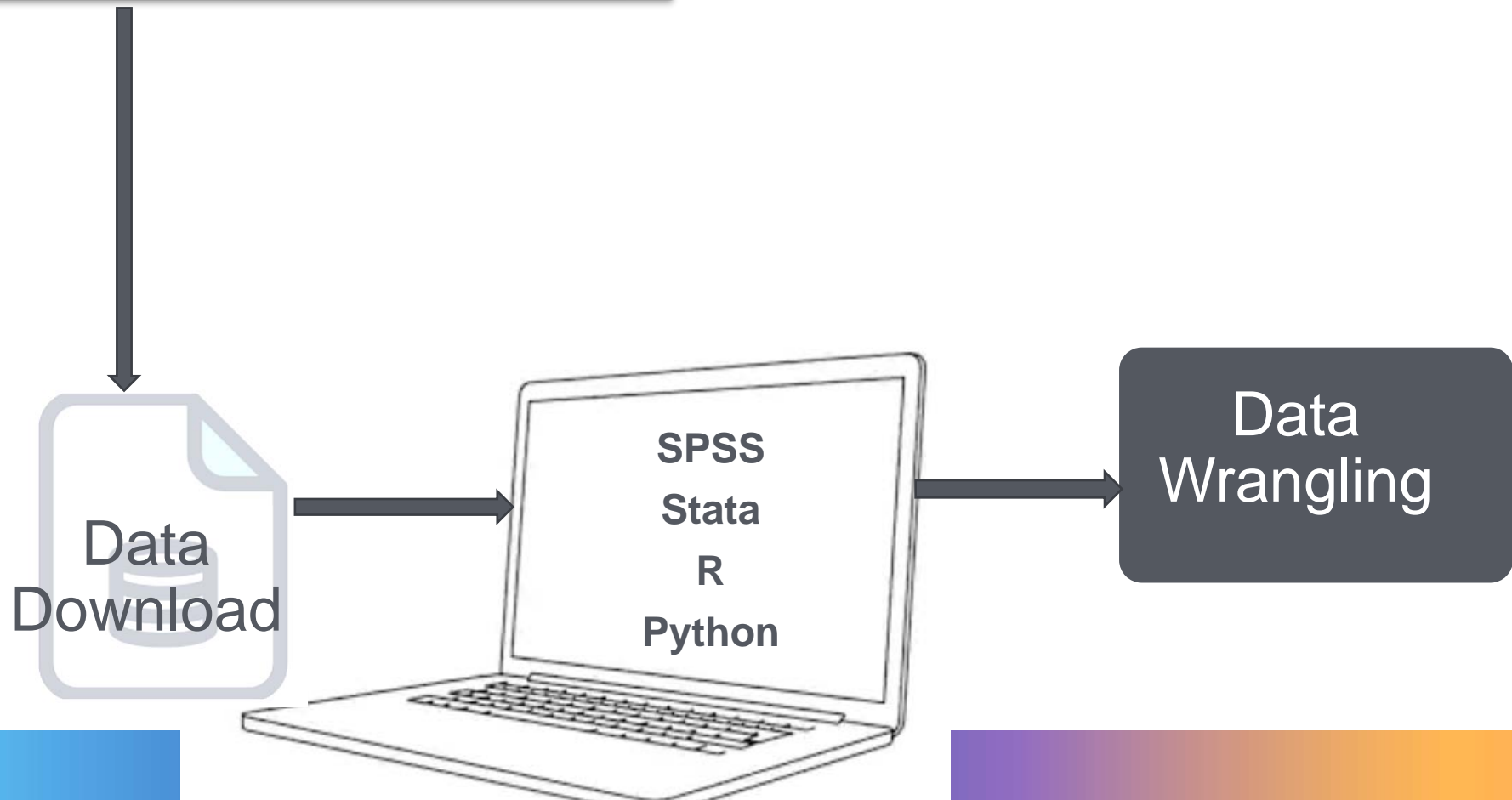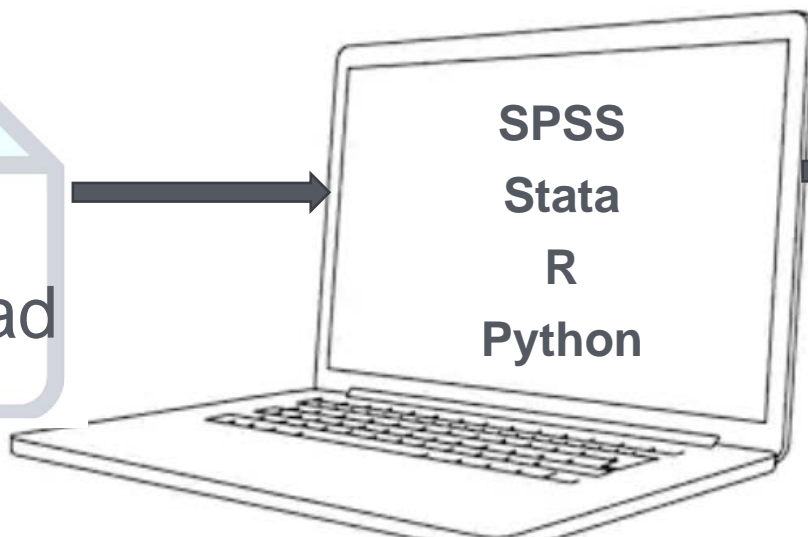
chapter1_2019.dat
chap1_2019.dat

Data Download

SPSS

Stata

R

Python

Data Wrangling

Data
Download

SPSS
Stata
R
Python

Data
Wrangling

Exploratory
Data
Analysis

Statistical
Modelling

Write-Up
Results

AMERICAN
SOCIOLOGICAL
REVIEW

# The Workflow

Analysing data without a planned and organised workflow can be compared to drinking and driving. In both situations it doesn't matter how careful you are, it is still highly likely to end in a wreck!

(Philip Stark UC Berkeley)

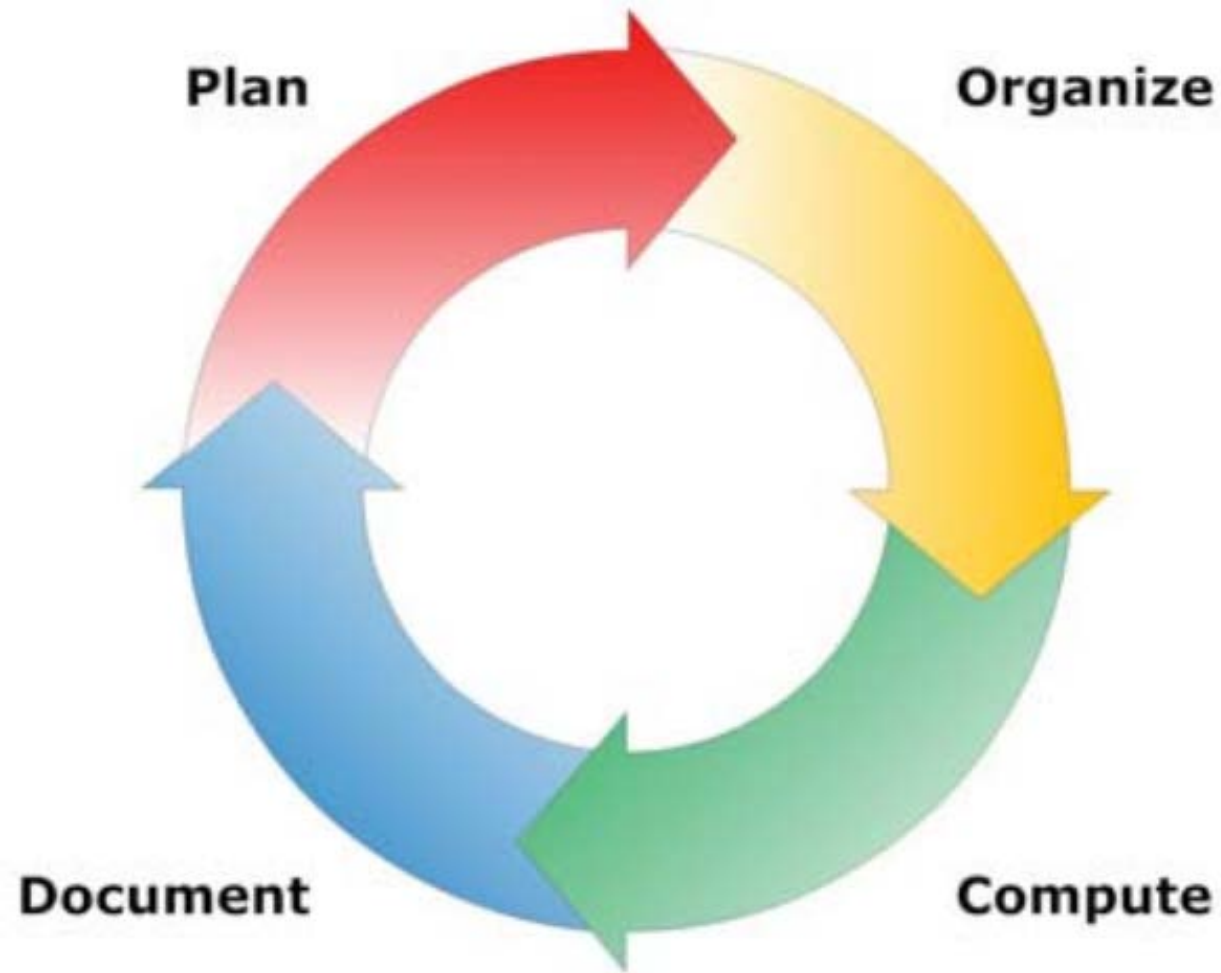Therefore just like drinking and driving, we strongly warn against not having a systematic workflow
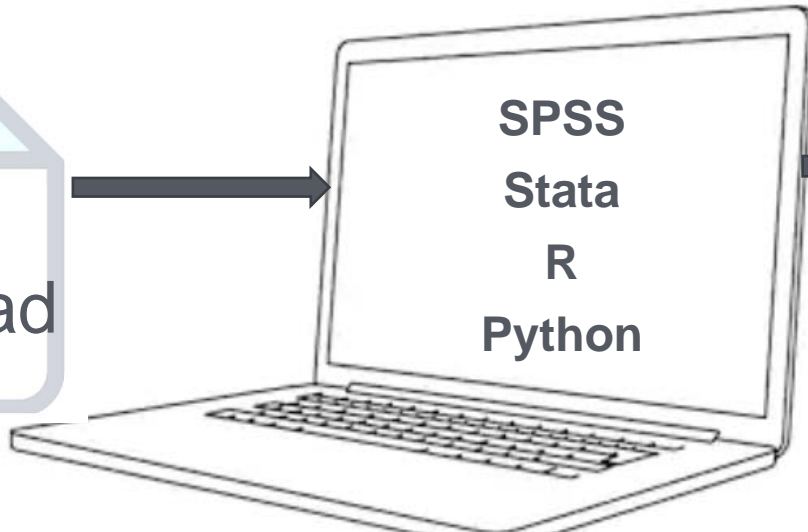
# The Workflow

Workflow should be planned and carefully orchestrated

Workflow MUST NOT be *adhoc*

(e.g. piece-meal, developed as a reaction to mistakes etc.)

# The Workflow Cycle

Data Download

SPSS

Stata

R

Python

Data Wrangling

# Drop down menus = no audit trail



GUIs will leave you in a sticky mess!

# Data Wrangling

- Selecting variables (surveys have large $k$)
- Operationalising measures (e.g. income; social class; education)
- Re-coding variables
- Selection cases (sub-samples in large surveys)
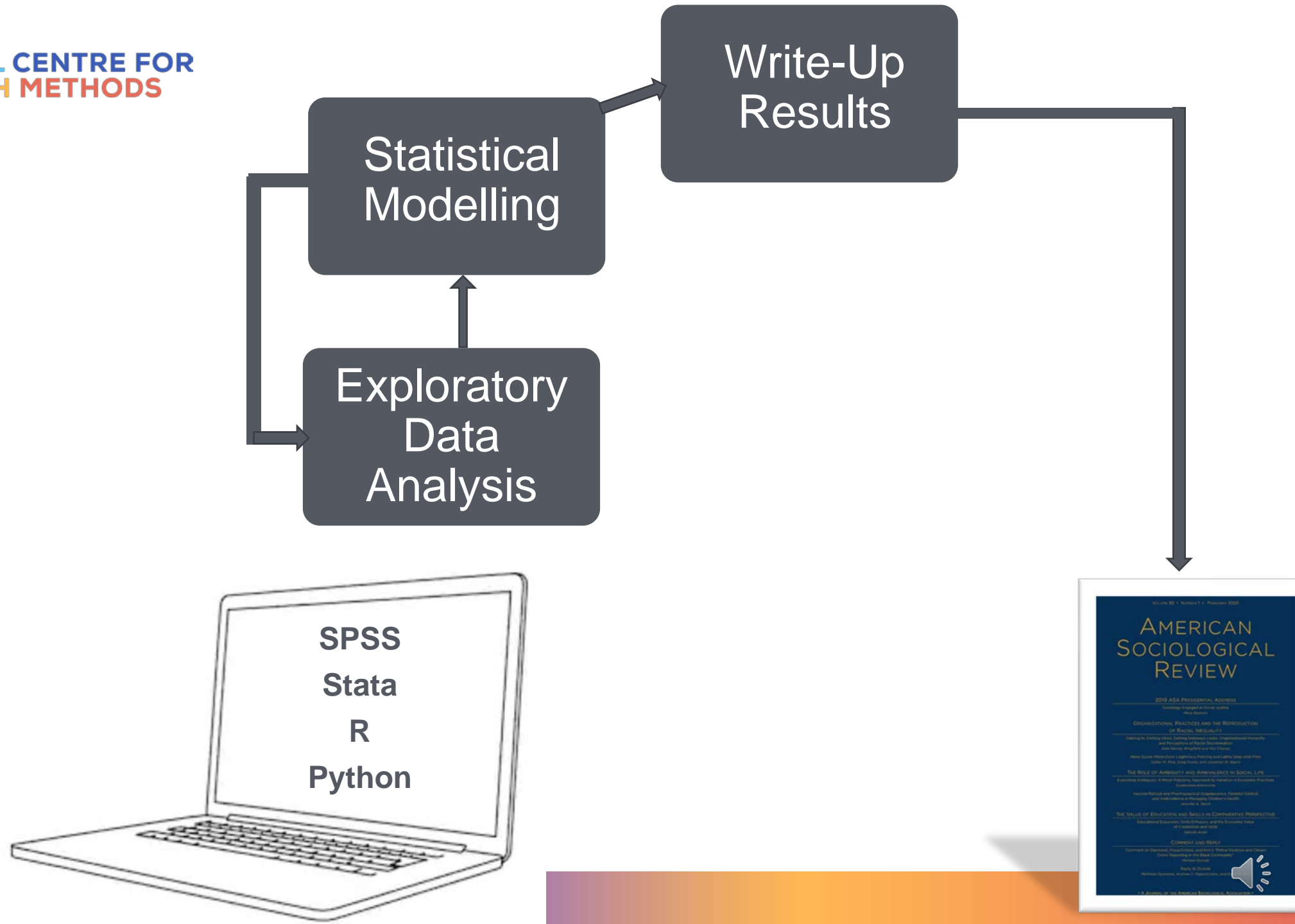- Missing data

# Data Wrangling

- Selecting variables (surveys have large $k$)

- Operationalising measures (e.g. money, social class, education)

- Recoding variables

- Selection cases (sub samples in large surveys)

- Missing data

**CAN THESE ACTIONS BE TRACED IN MY AUDIT TRAIL**

# Minor actions in the

## *Data Wrangling Phase*

# can have major consequences!

# Data Analysis

- Which cases
- Which variables (format)
- Missing data
- Estimation method (mle; gls)
- Weights and survey structure (svy)
- Setting a seed
- Number of quadrature points
- Which  version or library

# Data Analysis

- Which cases
- Which variables (format)
- Missing data
- Estimation method (mle, gls)
- Weights and survey structure (svy)
- Setting a seed
- Number of quadrature points
- Which version or library

**CAN THESE ACTIONS BE TRACED IN MY AUDIT TRAIL**

# Improving the Workflow

*Better supporting YOU and what YOU DO*

*Not changing you into something YOU ARE NOT*

# ALL SERIOUS WORK <u>MUST</u>

# BE REPRODUCIBLE!

# There <u>MUST</u> be an audit trail

# A Planned Workflow Has Benefits

# Four Pillars of Wisdom

- Accuracy
  - minimising information loss and errors in analyses and output

- Programming Efficiency
  - automation, maximising features in software

- Transparency
  - showing what you did, why, when, how

- Reproducibility
  - same results every time whoever or wherever
  - editing, rewriting reports or re-submission of papers

# Long's Law

*It is always easier to document today than it is tomorrow!*

*Corollary 1:*
*Nobody likes to write documentation*

*Corollary 2:*
*Nobody ever regrets having written documentation*

# Long's Law

*Has anyone in the history of data analysis ever said*

*"these files are too well documented"*

# Good News

- Improving the workflow with a modest amount of effort

- The less experience you have the better
  (start from the very beginning)

# Some Practical Strategies and Tactics

# Standard Directory Structure

# File Naming Protocols

File Name = name_date_depositor's initials_version_type

Therefore     **bhpsaindresp_20140506_vg_v1.dta**

Would be a

a. The British Household Panel Survey File "aindresp"
b. Deposited on 6th May 2014
c. Deposited by vg (Vernon Gayle)
d. Version v1
e. File type (e.g. a Stata .dta file)

# Analysis File Structure and Information

```
template* ×
 1    STOP
 2
 3    /**
 4
 5    ******************************************************************
 6
 7    Next Actions:
 8
 9
10
11
12    Author:
13
14
15    Project:
16
17
18    Sub-project:
19
20
21    Date of Next Meeting (or supervision):
22
23
24    Latest Update:
25
26
27    Previous Updates:
28
```

# Variable Naming Protocols

- Good data providers will tend to stick to clearly defined variable naming conventions

- Wave 1 of UKHLS *a_sex* is the gender variable, and in wave 2 of the study *b_sex* is the gender variable (UK Data Archive Study Number 6614)

- Youth Cohort Time Series for England, Wales and Scotland, 1984-2002 the variable *t0schtyp* is the type of school that the pupil attended in Year 11 (at time point t0 in the study) (UK Data Archive Study Number 5765)

# Variable Naming Protocols

1. Here are some things, both good and bad, which people have said about their 4th and 5th years at school. We would like to know what you think. Please tick a box for each one to say whether you agree or disagree.

Agree    Disagree

— School has helped to give me confidence to make decisions [1] [2]

By contrast in the less well curated Youth Cohort Study of England and Wales Cohort 4 (UK Data Archive Study Number 3107) this is the opaque variable *dx11_a*

# Estimating Work Time…

# Benefits of a Good Workflow

- Limits duplication of effort

- Systematic work minimises errors

- Helps to detect errors

- Editing work, resubmitting papers etc.

- Aids collaboration

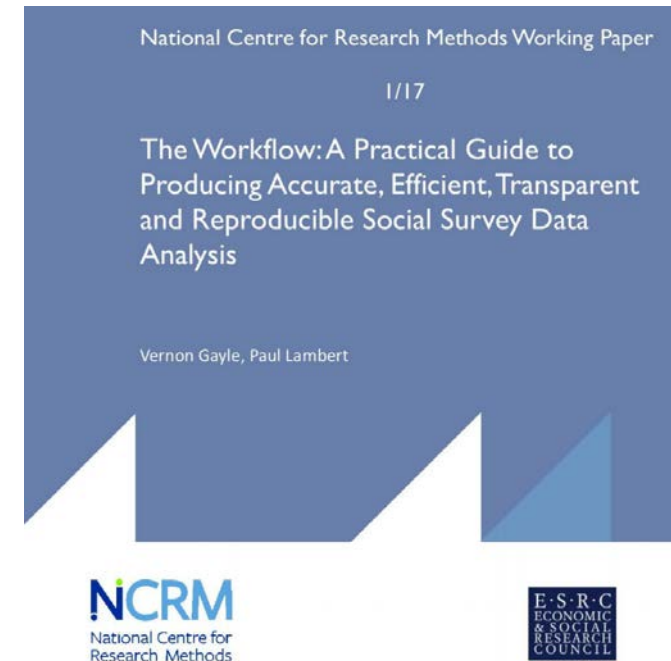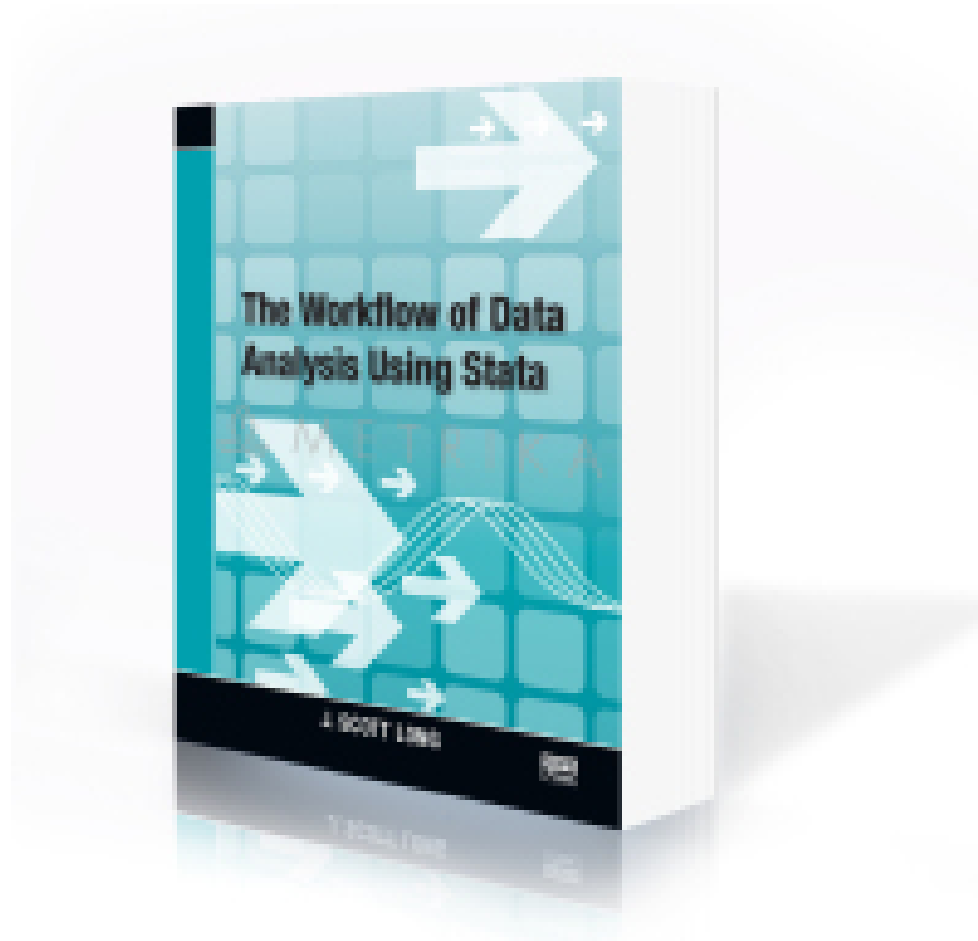- Supports additional (e.g. future) work

# What To Do…

- Audit trail

- Don't use drop menus ( a gui = a sticky mess)

- Systematic directory structure, file names, versions, variable names etc.

- Plan! don't do *ad hoc* work or work on the fly

- Write comments

National Centre for Research Methods Working Paper

1/17

The Workflow: A Practical Guide to Producing Accurate, Efficient, Transparent and Reproducible Social Survey Data Analysis

Vernon Gayle, Paul Lambert

http://eprints.ncrm.ac.uk/4000/

J. Scott Long has posted a really good pdf version of a talk on the workflow   http://www.ihrp.uic.edu/files/Workflow%20Slides%20JSLong%20110410.pdf

# How to cite this video

Gayle, V. (2020) *The Data Analysis Workflow.* Available at: https://www.ncrm.ac.uk (Accessed: day month year)