# Updating the Classics: a New Life for Old Data

Sharon Bolton

Data Publishing and Curation Manager

IASSIST & CARTO 2018: Once Upon a Data Point: Sustaining our Data Storytellers

Montreal, Canada, 31st May 2018

UK Data Service

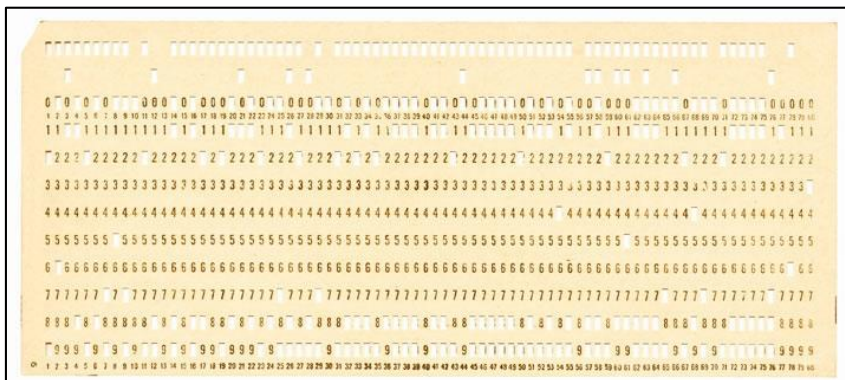UK • DATA ARCHIVE

University of Essex

# Setting the scene

- 2016 Brexit referendum result:
  - What had changed since 1975? New interest in revisiting old research

- Worked with researcher to find the data needed
  - Public attitudes to UK foreign policy in late 1960s/early 1970s
  - Lots of literature but less data
  - National Opinion Polls (NOP) data had most potential

- Roper Center holds extensive collection, behind paywall

- UK Data Service holds many – some not used for a long time
  - Paper documentation scanned to PDF and available on web, but
  - Data in older 'column binary' format, difficult to use

UK Data Service

# What is column binary data?

- Raw data once stored on computer punch cards - standard 80 columns of data occupies the 12 rows of each card. Data stored in this way are called column binary or multi-punch data, and allow more than one variable (in this case, survey question) to be stored in the same column  (adapted from Landman, 1996)



- Read by card reader machine that creates digital column binary data files
- Data can't be read by current statistical packages without conversion

# Column binary layout

- 12 columns → 80 rows ↓
- Usually, & = value 11, - = 12, 0 = 10, 1-9 = 4-12, Blank=missing
- But not always – check the data!

```
:***:***********************************************************************:***:
:Col:     &      -      0      1      2      3      4      5      6      7      8      9  Blank  :Col: Punches
:***:***********************************************************************:***:
:  1:                  911                                                                       :  1:    911
:  2:                         911                                                                :  2:    911
:  3:                  911                                                                       :  3:    911
:  4:                  271    288    352                                                         :  4:    911
:  5:                   65     73     62     97    122    107    108    111     81     85         :  5:    911
:  6:                   87     74     77     84     98    107     91    109     77    107         :  6:    911
:  7:                  105    100     85     89     98     90     86     87     91     80         :  7:    911
:  8:                  151    180    130    129    192    108     21                              :  8:    911
:  9:                   75     74    110    149     97     83     88     61     59    115         :  9:    911
: 10:                   54     82     87     97     52    141     72    140     86    100         : 10:    911
: 11:     82    133     50     72     91    106     62     75     19    108     58     55         : 11:    911
: 12:     57     18      7    223    149    102    144     48     54     39     47     23         : 12:    911
: 13:    911                   48                  261           108    266     64    164         : 13:   1822 Multipunched
: 14:      1                  199    552    156    290    535     83    291    398    219         : 14:   2724 Multipunched
: 15:    108      5      2    388    218     77     24     11    191     43     19     15      1  : 15:   1101 Multipunched
: 16:                         237    289    331     44      7                         1      4  : 16:    909 Multipunched
: 17:                    1    206    437    213             5      2      1      4      1     60  : 17:    870 Multipunched
: 18:                  323    163    258    513    153    235    317    284    343    322    130  : 18:   2911 Multipunched
: 19:                  254    148    206    179    202    202    232    240    245    267    166  : 19:   2175 Multipunched
: 20:                  140    215    201     92    207    195    162    177    137    132    252  : 20:   1658 Multipunched
: 21:                  162    354    226    102    314    247    162    181    151    147    220  : 21:   2046 Multipunched
: 22:                    1    298    229     35     22    314             3      2            14  : 22:    904 Multipunched
: 23:                         356    358    187      2      1      2      2      2      1      9  : 23:    911 Multipunched
: 24:                         224    537    130      3      3      1      1      2      1     21  : 24:    902 Multipunched
: 25:                         669    567    555    611    774    398    408    600    256     45  : 25:   4838 Multipunched
: 26:                         123    177    122    153     39    200    253     95    360    280  : 26:   1522 Multipunched
```

UK Data Service

# How did we get here?

- UK Data Archive 50 years old – large collection over lifetime

- Collection management requires time and resources, funding constraints had meant other work prioritised

- 'Data archaeologist' used documents from long-finished Gallup Poll project to develop conversion script for column binary > SPSS (Landman, 1996)

- Had used method before but never on this scale: over 50 datasets in NOP series, researcher needed c.30 in a timely fashion

UK Data Service

# Serendipity

- One-time extra funding agreed for collection management

- The mission – search and rescue
  - Find and convert the NOP series

- The quest
  - Recruited a curator and trained him in the tools to do the job
  - Background in data analysis
  - Programming in SPSS, using and manipulating syntax, running scripts

- The goal
  - Make the NOP series data available and easy to use

UK Data Service

# How do we do it?

- Make a map - tell the software where to look in the column binary data file and what to do with the information it finds

- Making the map depends on good metadata - if the map is correct you can find the treasure

- Did we have enough information to find our way? Write the script and see

- Documentation = curation notes, questionnaires, reports

UK Data Service

# Metadata dream

| CARD/COLUMN | VAR NO | TITLE | CODES |
|---|---|---|---|
| 1/1–4 | | Case Number | 0000 – 9999 |
| 1/5–6 | | Card Number | 01 Card 1 |
| 1/7–10 | | NOP Number | 6728 |
| 1/11 | 1 | Sex | 1 Male<br>2 Female, housewife<br>3 Female, non–housewife |
| 1/12 | 2 | Marital status | 1 Married<br>2 Single/widowed/divorced/separated |
| 1/13 | 3 | Head of household: Respondent | 1 Male head of household<br>2 Female head of household<br>3 Not head of household |
| 1/14 | 4 | Number of people in household | 1 1<br>2 2<br>3 3<br>4 4<br>5 5+ |
| 1/15 | 5 | Children | 1 Household has children under 16<br>2 Household has no children under 16 |
| 1/16 | 6 | Age | 1 16–20<br>2 21–24<br>3 25–34<br>4 35–44<br>5 45–54<br>6 55–64<br>7 65+ |
| 1/17 | 7 | Class | 1 A<br>2 B<br>3 C1<br>4 C2<br>5 DE |

*Table title: 074 NOP 6728 (10 – 15 APRIL 1973)*

# Metadata nightmare

# Drawing the map

- Write and format the script – no shortcuts
- Trial and error: run the script, check the results against the documentation – multipunch columns can cause errors
- Amend the script, run it again until data correct

# Completing the job

- Clean and label the data – recoding multi-punch variables, string characters to numeric, add metadata (variable and value labels)

- Enhance usability - apply robust UK Data Service curation standards

- Create preservation format (ASCII) and preservation metadata

- Create current standard dissemination formats – SPSS, Stata, tab-delimited text

- Upgrade scanned documentation – optical character recognition (OCR), PDF/A where possible

- Augment catalogue metadata

UK Data Service

# Before …

# After …

# End of Part One

- Special funding finished, back to conversion where and when we can

- Celebrate achievements so far
  - Impact! Researcher soon to publish book:

    Clements, B. *Public Opinion towards Foreign and Defence Policy in Britain, 1945-2017* (forthcoming, Routledge)

  - Not just NOP series, but >100 column binary datasets upgraded to current formats

  - Proven conversion methodology – scripts and algorithms work

  - Trained curator in useful data science skills

UK Data Service

# Looking forward to Part Two

- Remaining column binary datasets to convert

- Modify scripts to work in other software (R, SAS, others?)

- Make scripts available to others, user guide, GitHub?

- These data born digital to UK Data Archive, need machine to read hard copy cards

- Renewed interest in data rescue (Research Data Alliance (RDA) Data Rescue IG)

- New funding opportunities for collection management?

UK Data Service

# Questions

Sharon Bolton

[sharonb@essex.ac.uk](mailto:sharonb@essex.ac.uk)