Documenting and contextualising your data

Tom Ensom UK Data Archive

Looking after your research data, UK Data Archive 3-4 April 2014





Overview

A crucial part of making data user-friendly, shareable and with longlasting usability is to ensure they can be understood and interpreted by any user. This requires clear and detailed data description, annotation and contextual information.

Areas of coverage

- Why documentation is important
- Study-level documentation and context
- Data-level documentation
- Metadata
- Context debate



What is documentation?

- Data doesn't mean anything without documentation
 - A survey dataset becomes just a block of meaningless numbers
 - An interview becomes a block of contextless text
- Data documentation might include:
 - A survey questionnaire
 - An interview schedule
 - Records of interviewees and their demographic characteristics in a qualitative study
 - Variable labels in a table
 - Published articles that provides background information
 - Description of the methodology used to collect the data



Why document your data?

- Enables you to understand/interpret data when you return to it
- Needed to make data independently understandable i.e. reusable
- Helps avoid incorrect use/misinterpretation
- If using your data for the first time, what would a new user need to know to make sense of it?
- The UK Data Archive uses data documentation to:
 - supplement a data collection with documents such as a user guide(s) and data listing
 - ensure accurate processing and archiving
 - create a catalogue record for a published data collection



What should be captured?

Any useful documentation such as:

• final report, published reports, user guide, working paper, publications, lab books

Information on dataset structure

- Inventory of data files
- relationships between those files
- records, cases...

Variable-level documentation

- labels, codes, classifications
- missing values
- derivations and aggregations



What should be captured?

Contextual information about project and data

- background, project history, aims, objectives, hypotheses
- publications based on data collection

Data collection methodology and processes

- data collection process and sampling
- instruments used questionnaires, showcards, interview schedules
- temporal/geographic coverage
- data validation cleaning, error-checking
- compilation of derived variables
- weighting: factors and variables, weighting process
- secondary data sources used

Data confidentiality, access and use conditions

- anonymisation carried out
- consent conditions/procedures
- access or use conditions of data



Consider documentation early on

- Good data documentation and metadata depends on what you as the creator can provide
- Start gathering meaningful information from as early on in the research process as possible
- This consideration forms an important part of data management planning (which you will hear more on later in the course)



Quantitative study

- Smaller-scale study single user guide may contain compiled survey questionnaire, methodology information
- Example from Understanding Society, a bigger study many documents presented separately:

DOCUMENTATION		
Title	File Name	Size (KB)
Cognitive Ability Measures	6614_cognitive_ability_measures_v1-1.pdf	348
Revisions November 2013	6614_ukhls_2013_revisions.pdf	375
Wave 1 Adult Main Questionnaire	6614_understanding_society_wave1_questionnaire.v04.pdf	2802
Wave 2 Adult Main Questionnaire	6614_understanding_society_wave2_questionnaire_v04.pdf	3726
Waves 1-3 User Manual	6614_usermanual_wave1to3_v1-1.pdf	883
Wave 3 Youth Self-Completion Questionnaire (GB)	6614_w3_youthquestionnaire_gbritain_annotated.pdf	1469
Wave 1 Consent Package	6614_wave1_consent_package.pdf	645
Wave 1 Adult Self-Completion Questionnaire	6614_wave1_main_adult_sc_questionnaire.pdf	429
Wave 1 Youth Self-Completion Questionnaire	6614_wave1_main_youth_sc_questionnaire.pdf	750
Wave 1 Project Instructions for Interviewers	6614_wave1_project_instructions_interviewers.pdf	2426
Wave 1 Showcards	6614 wave1 showcards ndf	100



Qualitative study – user guide and doc

• A user guide could contain a variety of documents that provide context: interview schedule, transcription notes, even photos

<text><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></text>	QD1
<text><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></text>	ROTES ON THE INTERVIEW SCHEDULE
10. Proposed as a real of them able to recite the names of the children in the famile first places between them. If is useful in these places between them, the is a boot of the integration of the spaces between them, the is a boot of the integration of the spaces between them, the is a boot of the integration of the spaces between them. If is useful in the space is the space between them, the is a boot of the integration of the spaces between them. If is useful in the space is the space between them, and the space is	The bounded
Areas sometime if and it easier to write down or they much easier and makes of their statistic and doal at the present time. (a) Year respondents as not have the up of their father when the year born, and if they were born, and they were born, and they were born, and they were born, and were born, and they	1(c) wraphonetic are not often able to recite the names of the children in the family frequency purposed and the spaces hetween thes. It is areful in these earses to mak there through the space hetween the children of the space his and the spaces hetween the children who user older than him. Thus ask about the yronger ones. Keypondents are sometimes vague about the respective ages of tweir intervals. As, the comparison of the space of the space of the space of the space intervals. As, the comparison of the space of the space of the space of the space intervals. As, the space of the space intervals. As the space of the
(a) these respondences also have have the key of their failure when they were horm, and if they have how of her if they are hord they were horm, and if they have how of her is a set of the set	dents sometimes find it easier to write down or tell you the ages and names of their siblings, alive and dead, at the present time.
<section-header><section-header><section-header><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></list-item></section-header></section-header></section-header>	1(d) When respondents do not know the spc of their father when they were born, ask if they know how old their vas when he died (assuming he is dead) and what year that was. Or respondents the age their father was when he married and the date. Approximate dates will do.
 2. Densitie housing 3. Sector the house in which respondent spant the longest time is an ensemble for severy which is the intervention of the severy which is a lot of contact with the server and severy which is a lot of contact with the server and severy is a lot of the respondent of a lot of contact with the server and severy is a lot of the respondent of a lot of contact with the server and severy is a lot of the respondent of a lot of contact with the server and severy is a lot of contact with the server and severy is a lot of contact with the server and severy is a lot of contact with the server and severy is a lot of contact with the server and severy is a lot of contact with the server and severy is a lot of contact with the server and severy is a lot of contact with the server and severy is a lot of contact with the server and severy is a lot of contact with the server and severy is a lot of contact with the server and severy is a lot of contact with the server and severy is a lot of the server and sever and severy is a lot of the server and sever and sever and sever and sever server serve	1(e) See notes on 1(d).
 <u>A constrict fouring</u> <u>A constrict fouring</u><	
 4.9. sleet the house in which respondent spent the longest time he can remember above the large in a survive here. 9. Servants in this period who did net time in occen smally chromesen of owner hor survive here. The respondent search is the respondent. There were spendent to a solid local search is the local affect the respondent. The over survive here, the servant base the second search is a solid local search is a solid local search is a local search	Domestic Routine
 As remark in this period who did net line in norm smally characeness or women the same into de de rough (is is, to de the rough) humanost. There were about the second into the rough humanost. There were about the second is the de the vanhing and young girls who cans in to look after the roughed and young girls who cans in to look after the roughed and the second is the second is roughed and the roughed and the second is roughed and the second roughed and the s	2(a) Select the house in which respondent spent the longest time he can remember before leaving home.
 2(a) their existing a sensitive level of a far the younger children, took them out for sensitive sensitive	$2(\epsilon)$ Servants in this period who did not live in were usually charvonen or women who nome in "to do the rough housework. There were also washervonen who came in to do the rough housework. There were also children. Where the respondent as a child nome into a lot of contact with the servant particularly if the looked after the respondent, find out what the relationship was between their, the sort of things she did for the respondent, etc.
3. stall 3. stall 3. Some and women women torking day started early would often take momenting with the Mor Mor Markafast. Menn asking about seals find out den the regregedent took (and Morkafast), when asking about seals (at a weak asking to the second seals of the second secon	$2(\varsigma)$ Older children sometimes looked after the younger children, took them out for valks, saw them to school, etc
 M(c) were and weener winner devices in storted early would offen the concerting, with the offen the some time of the storted early week to be a some sort inclusion of the some particularly in class if and 2, it on angle-initial early of the some particularly in class if and 2, it on angle-initial early and and 2, it on angle-initial early of the some particularly in class if and 2, it on angle-initial early the some particularly in the some particular initial early and the source of the some particularly in the source of the some particular initial early and the source of the some particular initial early and the source of the so	3. Neals
3(h) Sometimes a person might take his plate and sit by the error of the fire durine a meat. Or a person in a hurry right snatch some food standing up.	McD from and downen working day started early would often take something, with them for brackfart. Backets the solar find out deep the respondent took food and wint be called those mealering about the solar too the respondent took food solar work to some, particularly is class 1 and 2, to an agricultural labourcer it is a mark content at about 11 a.m. Dimmer is the middly meal to the majority of respon- dents. To some, particularly is class 1 and 2 at about 2 rob p.m The too mark to be and the solar solar to the middly main to the middly solar to the and the solar meal to the day. To some, in class 1 and 2 at about 3 period a difference to the day. To some, in class 1 and 2 at about 3 period. If is an a aftermon to in that case. Supper may be a cup of cocean and some bread and cheese it any bas a need of two compress either how to the last gaal at about 3 period. If any of the solar to the solar t
	3(k) Sometimes a person might take his plate and sit by the corner of the fire during a meal. or a person in a burry right snatch some food standing up.





Qualitative study – data listing

• Data listing provides an at-a-glance summary of interview sets

Study Number 5407 Health and Social Consequences of the Foot and Mouth Disease Epidemic in North Cumbria, 2001 Mort, M.

The panel respondents for the study were divided into six population groups. The data list for the diary and interviews has been colour-coded accordingly for clarity, using the depositor's original colours:

		Group 3:			
	Group 2: Rural	Agricultural related	Group 4: Frontline		Group 6: Animal / Human
Group 1: Farmers	Business	occupations	Workers	Group 5: Community	Health Professionals

1. Interviews

Respondent ID	Population Group	Date of Birth	Gender	Occupation	Interview summary	Place of Interview
					Family and	
					background,career and work,	
	Group 6: Animal /				arrangements during FMD	
	Human Health				epidemic and perceptions of	North Cumbria, respo
PM02	Professionals	1975	M	Veterinary Surgeon	situation	home
					Family and	
					background,career and work,	
	Group 6: Animal /				arrangements during FMD	
	Human Health				epidemic and perceptions of	
PM03	Professionals	1966	F	Veterinary Surgeon	situation	North Cumbria
					Family and	
					background,career and work,	
	Group 6: Animal /				arrangements during FMD	
	Human Health				epidemic and perceptions of	North Cumbria, respo
PM07	Professionals	1964	F	Veterinary practice manager	situation	home
					Family and	
					terration and the second second	



Data-level documentation

- Certain types of data file may contain important information which should be preserved:
 - variable/value labels; document metadata; table relationships and queries in relational databases; GIS data layers/tables
- Some examples:
 - SPSS: variable attributes documented in Variable View (label, code, data type, missing values)
 - MS Access: relationships between tables
 - ArcGIS: shapefiles (layers) and tables in geodatabase; metadata created in ArcCatalog
 - MS Excel: document properties, worksheet labels (where multiple)



Embedded data-level metadata in SPSS file

🔝 hse09ai.	.sav [DataSet2] -	PASW Statistics	Data Editor			
File Edit	<u>V</u> iew <u>D</u> ata	Transform An	alyze Direc	t Marketing	Graphs Utilities Add-ons Window Help	
		5 7				
	Name	Туре	Width	Decimals	Label	Values Missing
175	quala10	Numeric	2	0	Which of the qualifications on this card do you have? 10	{-9, No ans991
176	activb	Numeric	2	0	Activity status for last week	{-9, No ans991
177	empstat	Numeric	2	0	Manager/Foreman	{-9, No ans991
178	everjob	Numeric	2	0	Ever had paid employment or self-employed	{-9, No ans991
179	ftptime	Numeric	2	0	Full-time or part-time	{-9, No ans991
180	howlong	Numeric	2	0	How long have you been looking	{-9, No ans991
181	wkstrt2	Numeric	2	0	Able to start work within 2 weeks (Government training scheme)	{-9, No ans991
182	wklook4	Numeric	2	0	Looking paid work/govt scheme last 4 weeks	{-9, No ans991
183	nemplee	Numeric	2	0	Number employed at place of work	{-9, No ans991
184	nssec	Numeric	5	1	NS-SEC - long version (harmonised)	{-9.0, No a99.01.0
185	othpaid	Numeric	2	0	Ever had other employment (waiting to start work)	{-9, No ans991
186	payage	Numeric	3	0	Age when last had a paid job	{-9, No ans991
187	paylast	Numeric	4	0	Year left last paid job	{-9, No ans991
188	paymon	Numeric	2	0	Month last left paid job	{-9, No ans991
189	sclass	Numeric	2	0	Social Class	{-9, No ans991
190	seg	Numeric	2	0	Socio-Economic Group	{-9, No ans991
191	snemplee	Numeric	2	0	Self employed, how many employees	{-9, No ans991
192	age	Numeric	3	0	Age last birthday	{-9, No ans991
	1					1
Data View	Variable View	/				



ce

Data-level documentation: variable names

- All structured, tabular data should have cases/records and variables adequately documented with names, labels and descriptions
- Variable names might include:
 - question number system related to questions in a survey/questionnaire e.g. Q1a, Q1b, Q2, Q3a
 - numerical order system

e.g. V1, V2, V3

 meaningful abbreviations or combinations of abbreviations referring to meaning of the variable

e.g. oz%=percentage ozone, GOR=Government Office Region, moocc=mother occupation, faocc=father occupation

 for interoperability across platforms - variable names should be max 8 characters and without spaces



Data-level documentation: variable labels

- Similar principles for variable labels:
 - be brief, max. 80 characters
 - include unit of measurement where applicable
 - reference the question number of a survey or questionnaire

e.g. variable 'q11hexw' with label 'Q11: hours spent taking physical exercise in a typical week' - the label gives the unit of measurement and a reference to the question number (Q11b)

- Codes of, and reasons for, missing data
 - avoid blanks, system-missing or '0' values
 - e.g. '99=not recorded', '98=not provided (no answer)', '97=not applicable', '96=not known', '95=error'
- Coding or classification schemes used, with a bibliographic ref
 - e.g. Standard Occupational Classification 2000 a list of codes to classify respondents' jobs; ISO 3166 alpha-2 country codes - an international standard of 2-letter country codes



Data-level documentation: transcripts

- Qualitative data/text documents:
 - interview transcript speech demarcation (speaker tags)
 - document header with brief details of interview date, place, interviewer name, interviewee details, context









"created_at": "Mon Jun 10 21:09:19 +0000 2013", "id": 344199622916448260, "id_str": "344199622916448256", "text": "Need to catch up? Our complete #NSAFiles coverage is here: http://t.co/iZPkknopxk", "source": "SocialFlow", "truncated": false, "user": { "id": 16042794, "id_str": "16042794" "name": "GuardianUS", "screen_name": "GuardianUS", "location": "New York", "description": "Featuring the Guardian's US coverage, conversations and reporters.", "url": "http://t.co/eqPiqNUSme", "protected": false, "followers_count": 55597, "friends_count": 509, "listed_count": 2414. "created_at": "Fri Aug 29 14:52:08 +0000 2008" "favourites_count": 860, "utc_offset": -18000. "time_zone": "Eastern Time (US & Canada)", "geo_enabled": true, "verified": true, "statuses_count": 41567, "lang": "en", 'contributors_enabled": false, "is_translator": false.



```
<dataKind>Semi-structured diaries</dataKind>
  </sumDscr>
 </stdyInfo>
▼<method>
 <dataColl>
    <timeMeth>Cross-sectional (one-time) study</timeMeth>
    <sampProc>Volunteer sample</sampProc>
   ▼<sampProc>
      An independent professional recruited respondents to a
      demographic profile agreed by the project steering group. See
      documentation for further details.
    </sampProc>
   ▼<deviat>
      42 individual interview transcripts, 40 diaries, 6 focus group
      transcripts and 1 audiomontage transcript. </br>The collection
      also includes 42 individual interview audio files, 7 focus group
      audio files, 1 audiomontage and 7 newsletters, but access to
      these is subject to permission from the depositor.
    </deviat>
   ▼<collMode>
      Face-to-face interview; Diaries; Compilation or synthesis of
      existing material; Focus group; Audio recording
    </collMode>
    <sources/>
    <weight>Not applicable</weight>
    <cleanOps>A</cleanOps>
  </dataColl>
 </method>
```

Catalogue

UK Data Service data catalogue record for:

Health and Social Consequences of the Foot and Mouth Disease Epidemic in North Cumbria, 2001-2003

Documentation Relate	ed Studies Publications 🐨 Download/Order Get full DDI XML
TITLE DETAILS	^
SN:	5407
Title:	Health and Social Consequences of the Foot and Mouth Disease Epidemic in North Cumbria, 2001-2003
Alternative title:	Health and Social Consequences of the 2001 Foot and Mouth Disease Epidemic
Persistent identifier:	10.5255/UKDA-SN-5407-1
Depositor:	Mort, M., Lancaster University. Institute for Health Research
Principal investigator(s):	Mort, M., Lancaster University. Institute for Health Research
Sponsor(s):	Department of Health
Grant number:	121/7499

SUBJECT CATEGORIES

Community and urban studies - Society and culture Rural life - Society and culture ~



Metadata – data about data

- In some ways, just another kind of documentation
- But much more highly structured
- Standard data collection metadata includes:
 - Components of a bibliographic reference
 - Core information that a search engine indexes to make the data findable
- International standards/schemes
 - Data Documentation Initiative (DDI)
 - ISO19115
 - Dublin Core
 - Metadata Encoding and Transmission Standard (METS)
 - Preservation Metadata Maintenance Activity (PREMIS)



Where to begin?!





Metadata at the UK Data Archive

- Metadata for archived datasets at should include:
 - Core fields: title, abstract, details of data owner/creator
 - Administrative: Funding information source and award number, copyright holder
 - Detailed descriptive info: temporal coverage (data collection start and end dates), geographic coverage (country, region, longitude/latitude), keywords and subject categories
 - Methodological: sample size/units, methodology
 - Data availability/access conditions
 - Publications/references
 - Digital Object Identifier (DOI)
- Created from data deposit form/tool and information/documentation submitted by data owners/researchers
- UK Data Service: DDI metadata, rich detailed content

ukdataservice.ac.uk/manage-data/document/metadata.aspx



How to create metadata for data

- Can be compiled using data deposit forms/tools
- Currently not many available that are user friendly and maintained
- May be better to take things into your own hands create a spreadsheet!
- Data Documentation Initiative (DDI) documentation can be created in software packages using certain DDI tools: tools.ddialliance.org
 - Colectica Designer for survey data <u>www.colectica.com/software/designer</u>
 - Convert SPSS internal metadata to DDI using Nesstar Publisher www.nesstar.com/software/publisher.html
 - German Institute for Educational Progress (IQB) educational data codebooks <u>www.iza.org</u>



Metadata entry for UK Data Service ReShare

	Logged in as Thomas Ensom Logout UK Data Service home Help About Contact
JK Data Service ReShare	Home Legal
	Edit collection: Data Collection #851298
ly data	
lanage records	Terms and conditions - Award details - People - Data collection - Upload - Deposit
Profile	
leview	* Data collection title ?
dmin	
dit page phrases	
	+ Alternative title
	* Data collection description ?



Exercise – how to document this data?

You carry out research on the public understanding of climate change and associated risks in the UK. Your data-generating research consists of:

- An online survey with 2000 invited members of the public in the UK to assess their understanding of climate change and climate change risks, as well as their sources of information.
- Interviews with 20 key stakeholders in climate policy and science communication.
- Qualitative content analysis of secondary data taken from newspapers and popular science journals, evaluating reporting about climate change in the media.

Data resulting from the online survey are transferred to SPSS for analysis.

Interviews are audio-recorded and transcribed into MS Word. Transcripts are imported into the NVivo software for content analysis.

Secondary textual data from newspapers and journals are also imported into NVivo for content analysis.

How would you document your resulting research data, to enable their future use by other researchers?



Exercise – some possible answers

For the survey data file, ensure that the SPSS file contains the full question text as variable label for each corresponding variable, or as detailed a description of the question text as possible. Variable names should consist of meaningful codes. Variable attributes are clearly defined, complete and without any abbreviations, explaining codes, categories and missing data values for each variable.

The online survey questionnaire is exported as a PDF file to complement the SPSS data file.

Each interview transcript contains an introductory paragraph providing the context and setting for the interview. The collection of 20 transcripts is accompanied by a data listing. Alternatively in NVivo a classification is created for all interviewees and interviews, capturing: for interviewees', relevant demographic and background characteristics (identifier or pseudonym) such as gender, age, profession, organisation and communication medium used; and for interviews, date, place and interviewer name.

A list or table of bibliographic references contains all sources of information used for the secondary analysis of media content.

A published article provides background information on the research methods used, sampling and so on.



Questions?

Contact details:

Tom Ensom UK Data Archive

datasharing@ukdataservice.ac.uk ukdataservice.ac.uk/manage-data.aspx

