

Documentation Quality Assessment in Ex-Post Survey Harmonization: Implications for Comparative Research

Ilona Wysmulek

ilona.wysmulek@ifispan.waw.pl

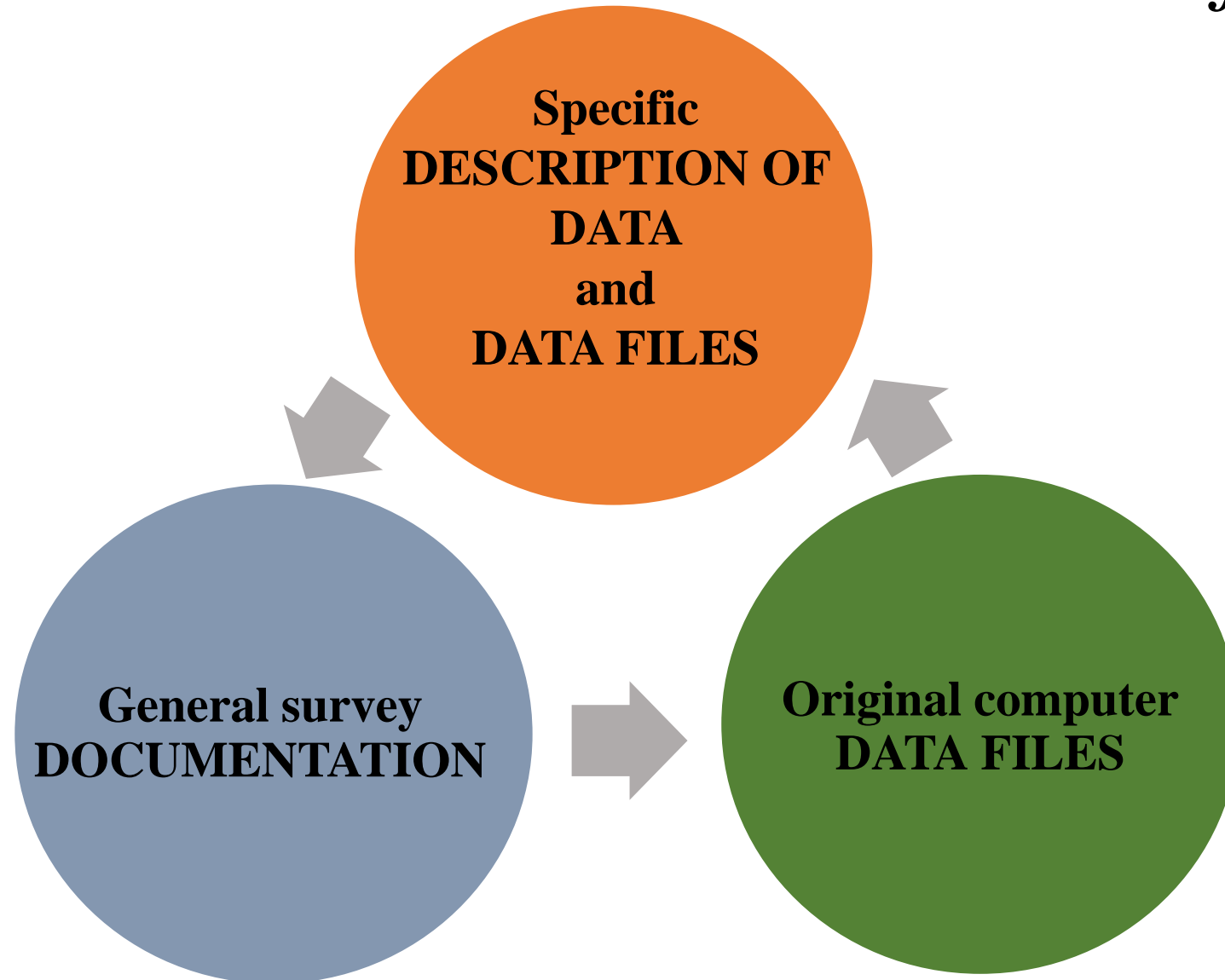
Olena Oleksiyenko

olena.oleksiyenko@gmail.com

Polish Academy of Science, Institute of Philosophy and Sociology



Quality Realms in the *Data Harmonization Project*



Working with: data AND documentation

- Interpretability of the data: does good documentation mean good quality of data?
- “coding, editing, ... and other data processing activities that follow the data collection phase” might be the source of errors (Groves 1989:12). [recognized in the literature]
- But... „[processing errors are] too rarely included in models of survey error” (Groves 2010: 869)

Sample

#	Target Variables	# Source var. /wave		# Waves
		Min	Max	
1	Gender	1	1	89
2	Age	1	3	89
3	Birth Year	0	1	30
4	Education levels	0	18	76
5	Schooling years	0	2	75
6	Trust in parliament	0	1	77
7	Participation in demonstration	0	4	75

Sample: 687 source variables matched to target variables

Sources

Data gathered from:

- a) codebook and/or questionnaire
- b) SPSS dictionary and original ,raw' data

Basic information: a) variable name b) exact question formulation c) variable label (codebook/SPSS dictionary) d) value labels (codebook/questionnaire/SPSS dictionary) e) exact values in the data

Data and Documentation Discrepancies: 8 types

VARIABLE LABEL	VARIABLE VALUE	INFORMATION
1.Misleading Label	2. Illegitimate Value	7.Insufficient Information
	3.Misleading Value	8.Translation Issues
	4.Value Discrepancy	
	5.Contradictory Values	
	6.Lack of Labels	

Example:

Survey	Question	Codebook	SPSS dictionary	DATA
ASB/2	Have you attended a demonstration or protest march?	Once 1 More than once 2 Never Done 3 CC 8 DA 9	1 Once 2 More than once 3 Never 8 Can't choose 9 Decline to answer	null 1 2 3 7 8 9

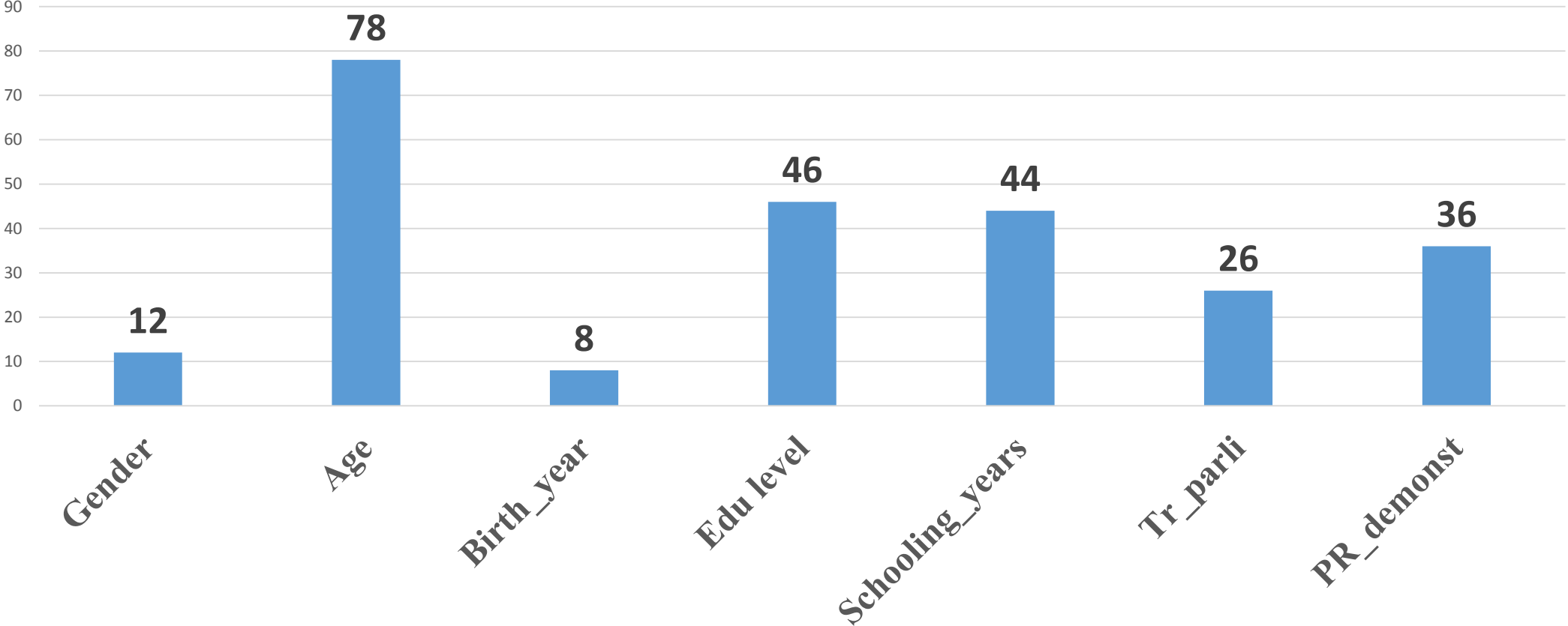
Frequencies of Discrepancies

[total: 687 source variables]

		<i>Frequency</i>
Total number of variable containing errors		197
Number of errors	1	146
	2	49
	3	2
Total number of errors		250

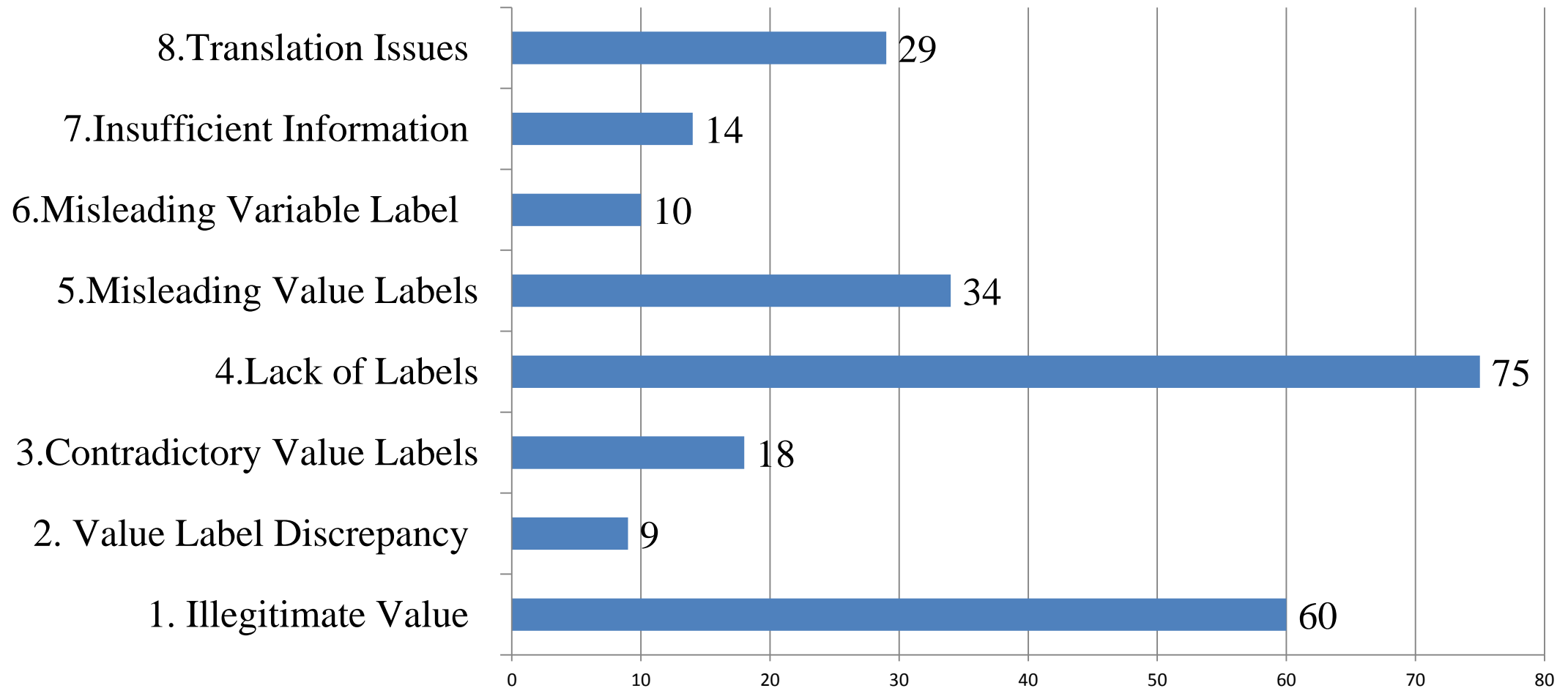
Discrepancies per Target Variables

[total =250 discrepancies]

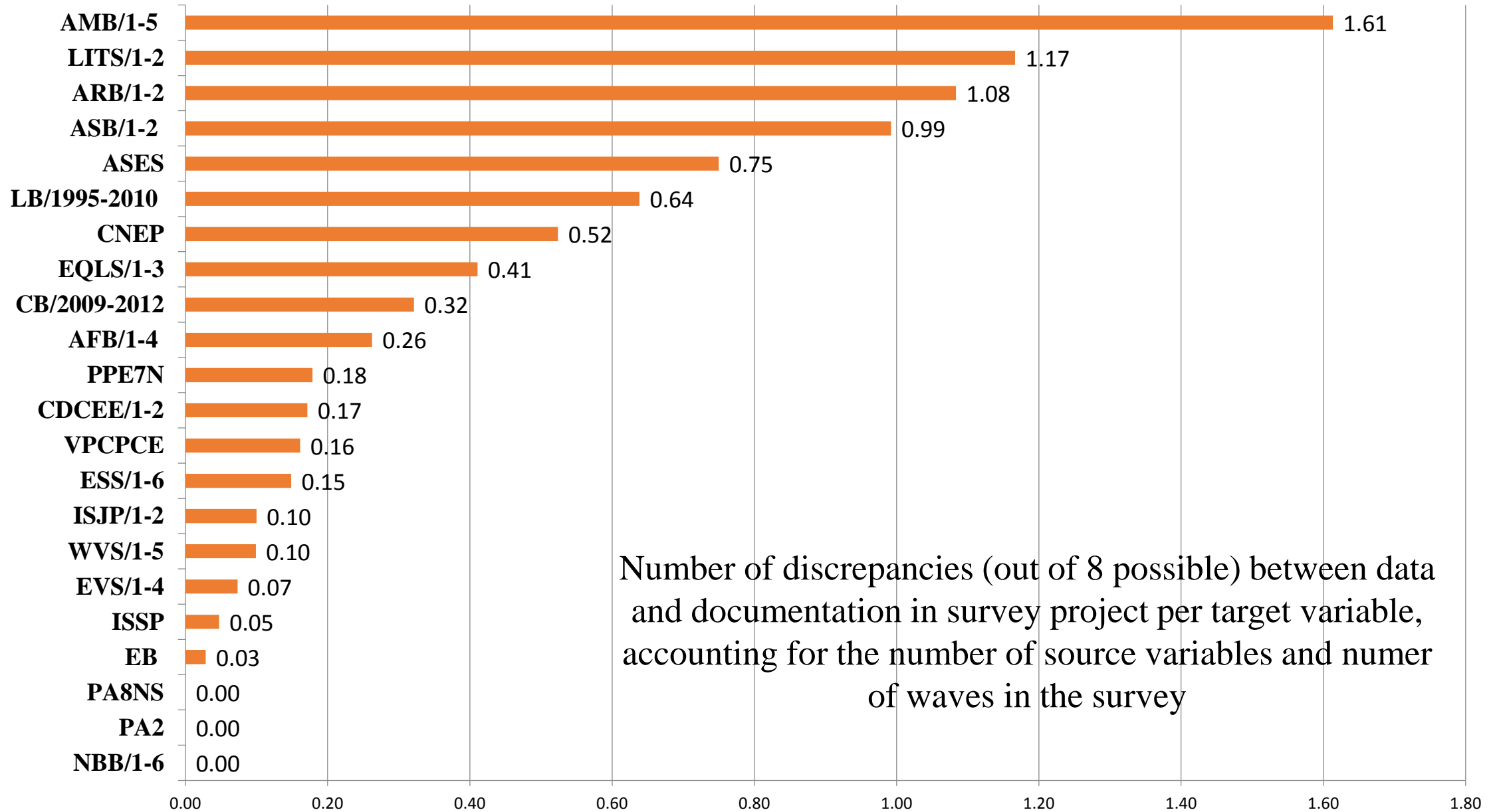


Types of Discrepancies

[total =250 discrepancies]



QUALITY INDEX per Survey Project



Discussion

- Discrepancies between data and documentation = decrease of interpretability of data
- How discrepancy typology can be used?
- Potential quality threat: different weight of types of discrepancies

- 20% of discrepancies out of all variables checked – is it a lot?