

VU Research Portal

CEDAR: The Dutch Historical Censuses as Linked Open Data

Meroño-Peñuela, Albert; Ashkpour, Ashkan; Guéret, Christophe; Schlobach, Stefan

published in
Semantic Web
2017

DOI (link to publisher)
[10.3233/SW-160233](https://doi.org/10.3233/SW-160233)

document version
Publisher's PDF, also known as Version of record

document license
Article 25fa Dutch Copyright Act

[Link to publication in VU Research Portal](#)

citation for published version (APA)

Meroño-Peñuela, A., Ashkpour, A., Guéret, C., & Schlobach, S. (2017). CEDAR: The Dutch Historical Censuses as Linked Open Data. *Semantic Web*, 8(2), 297-310. <https://doi.org/10.3233/SW-160233>

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

E-mail address:
vuresearchportal.ub@vu.nl

CEDAR: The Dutch historical censuses as Linked Open Data

Editor(s): Pascal Hitzler, Wright State University, USA

Solicited review(s): Ziqi Zhang, University of Sheffield, UK; Eetu Mäkelä, Aalto University, Finland & University of Helsinki, Finland & University of Oxford, UK; one anonymous reviewer

Albert Meroño-Peñuela^{a,b,*}, Ashkan Ashkpour^c, Christophe Guéret^{a,b} and Stefan Schlobach^a

^a *Department of Computer Science, VU University Amsterdam, De Boelelaan 1081a, 1081HV Amsterdam, The Netherlands*

E-mails: albert.merono@vu.nl, c.d.m.gueret@vu.nl, k.s.schlobach@vu.nl

^b *Data Archiving and Networked Services, Anna van Saksenlaan 10, 2593HT Den Haag, The Netherlands*

E-mails: albert.merono@dans.knaw.nl, christophe.gueret@dans.knaw.nl

^c *International Institute of Social History, Cruquiusweg 31, 1019AT Amsterdam, The Netherlands*

E-mail: ashkan.ashkpour@iisg.nl

Abstract. Here, we describe the CEDAR dataset, a five-star Linked Open Data representation of the Dutch historical censuses. These were conducted in the Netherlands once every 10 years from 1795 to 1971. We produce a linked dataset from a digitized sample of 2,288 tables. It contains more than 6.8 million statistical observations about the demography, labour and housing of Dutch society in the 18th, 19th and 20th centuries. The dataset is modeled using the RDF Data Cube, Open Annotation, and PROV vocabularies. These are used to represent the multidimensionality of the data, to express rules of data harmonization, and to keep track of the provenance of all data points and their transformations, respectively. We link observations within the dataset to well known standard classification systems in social history, such as the Historical International Standard Classification of Occupations (HISCO) and the Amsterdamse Code (AC). The three contributions of the dataset are (1) an easier access to integrated census data for historical researchers; (2) richer connections to related Linked Data resources; and (3) novel concept schemes of historical relevance, like classifications of historical religions and historical house types.

Keywords: Social history, census data, Linked Open Data, RDF Data Cube

1. Introduction

In this document we describe the CEDAR dataset, a five-star Linked Open Data conversion of the Dutch historical censuses dataset.¹

The Dutch historical censuses were collected from 1795 until 1971, in 17 different editions, once every 10 years. The government counted the entire population of the Netherlands, door by door, and aggregated the results in three different census types: demographic (age, gender, marital status, location, belief),

occupational (occupation, occupation segment, position within the occupation), and housing (ships, private houses, government buildings, occupied status). After 1971, censuses stopped from being collected from the entire population, mostly due to social opposition, and authorities switched to use municipal registers and sampling. Three facts make the 1795–1971 dataset self-contained and of special interest to historians and social scientists: it is based on counting the whole Dutch population, instead of sampling; it provides an unprecedented level of detail, hardly comparable to modern censuses due to privacy regulations; and the survey microdata from which the aggregations were originally built is almost entirely lost.

* Corresponding author. E-mail: albert.merono@vu.nl.

¹ See <http://www.volkstellingen.nl/>.

The census aggregations were written down in tables and published in books, archived by the Central Bureau of Statistics² (CBS) and the International Institute of Social History³ (IISH). In an effort to improve their systematic access, part of the tables in the historical censuses books have been digitized as images in several projects between the CBS, the IISH and several institutes of the Royal Netherlands Academy of Arts and Sciences⁴ (KNAW), such as Data Archiving and Networked Services⁵ (DANS) and the Netherlands Interdisciplinary Demographic Institute⁶ (NIDI). Beyond digitisation, these projects have translated part of this dataset, by manual input, into more structured formats. As a result, a subset of the dataset is available as a collection of 507 Excel spreadsheets, containing 2,288 census tables.

Challenges. The historical Dutch censuses have been collected for almost two centuries with different information needs at given times [1]. Census bureaus are notorious for changing the structure, classifications, variables and questions of the census in order to meet the information needs of a society. Not only do variables change in their semantics over time, but the classification systems in which they are organized also change significantly, making it extremely cumbersome to use the historical censuses for longitudinal analysis. The structures of the spreadsheets and changing characteristics of the census currently do not allow comparisons over time without extensive manual input of a domain expert. Even when converted into Web structured data, the need for harmonization across all years is a pre-requisite in order to enable greater use of the census by researchers and citizens.

Contributions. The goal of CEDAR⁷ is to integrate the Dutch historical censuses in these spreadsheets using Web technologies and standards; to publish the result of this integration as five-star Linked Open Data; and to investigate how semantic technologies can improve the research workflow of historians. Concretely, the main contributions of the dataset are:

- It is the first historical census data made available as LOD, integrated and Web-enabled from heterogeneous sources;

- it is linked to other datasets in the LOD cloud to improve its exposure and richness;
- it is released together with auxiliary resources, like historical classification schemes and integration mappings.

Additionally, the Dutch historical censuses Linked Open Data comes with the following features:

- Historical statistics on two centuries of Dutch history, fully compliant with RDF Data Cube [5];
- Standardization and harmonization procedures encoded using Open Annotations [22];
- Full tracking of provenance in all activities and consumed/produced entities as of PROV [9];
- Dereferenceable URIs;⁸
- A human browseable web front-end;⁹
- Dataset live statistics.¹⁰

The rest of the paper is organized as follows. In Section 2 we survey related work. In Section 3 we describe our conversion pipeline. In Section 4 we provide a full description of the data model and the use of established vocabularies, along with the quantity, quality and purpose of links to other datasets. In Section 5 we argue the importance of the dataset and its availability, including plans for long term preservation of the produced Linked Open Data. We discuss the five-star conformance of the dataset and its known shortcomings in Section 6.

2. Related work

Related work can be divided into (a) workflows and tools converting Excel/CSV data to RDF data; (b) other projects publishing statistics as RDF Data Cube; and (c) methods for enriching these tabular-converted RDF graphs.

There are many tools that convert tabular data to RDF.¹¹ CSV and HTML tables can be turned into RDF with dedicated tools [13,20]. Larger frameworks, like Open Refine + DERI's RDF plugin [7,19], Opencube [12] and Grafter¹² cover even more tabular and structured data formats, like Excel, JSON, XML, RDF, and Google Data documents. However, non of these are

²See <http://www.cbs.nl/>.

³See <http://www.iisg.nl/>.

⁴See <http://www.knaw.nl/>.

⁵See <http://www.dans.knaw.nl/>.

⁶See <http://www.nidi.knaw.nl/en/>.

⁷See <http://cedar-project.nl/> and <http://www.ehumanities.nl/>.

⁸See <http://lod.cedar-project.nl:8888/cedar/page/harmonised-data-dsd>.

⁹See <http://lod.cedar-project.nl/cedar/>.

¹⁰See <http://lod.cedar-project.nl/cedar/stats.html>.

¹¹<https://github.com/timrdf/csv2rdf4lod-automation/wiki>, <http://www.w3.org/wiki/ConverterToRdf>.

¹²See <http://grafter.org/>.

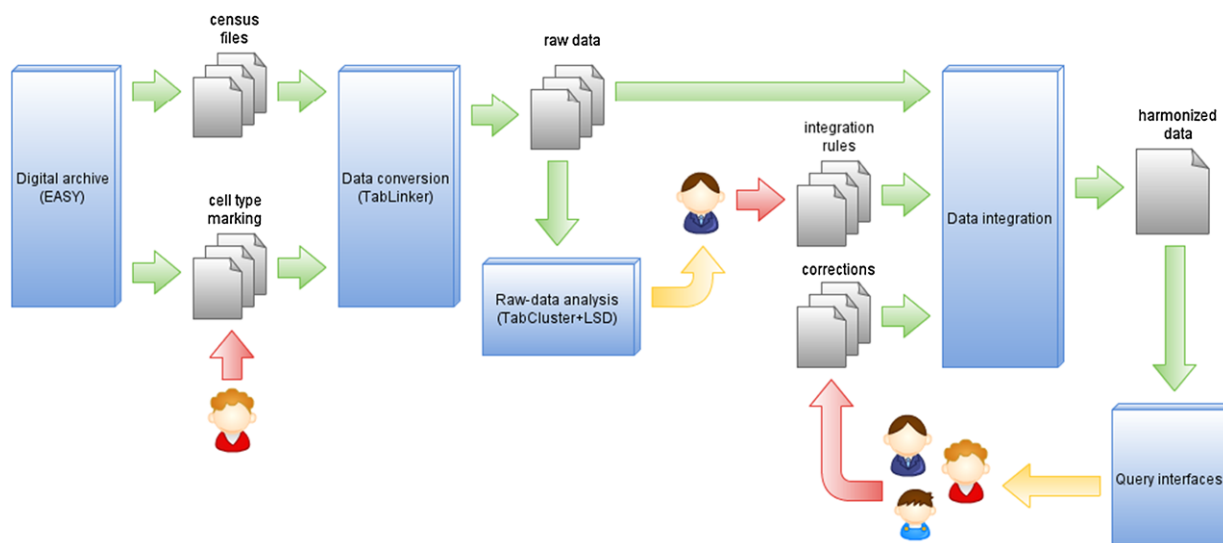


Fig. 1. Integration pipeline for the CEDAR data. The workflow starts at the archiving system, where the original Excel files are stored and retrieved using its API. Raw data is produced after interpreting complex table layout. These raw data are later transformed into harmonized data by applying integration rules encoded as Open Annotations. Red arrows indicate that manual input is required.

well suited for the conversion of historical tables. Data are sorted in these tables in a “eccentric” layout with spanning and hierarchical headers, multidimensional views and arbitrary data locations [1,17], which does not match the regularity of modern tables (i.e. one observation/record per row). To the best of our knowledge, only TabLinker [16] supports tables with these characteristics.

Other projects publish statistical datasets as RDF Data Cube. LSD Dimensions [14] provides a comprehensive index of statistical dimensions and Data Structure Definitions (DSDs) linking to those datasets. Remarkable ones include Linked Statistical Dataspaces [4] (with World Bank, European Central Bank and UNESCO Institute for Statistics data, among others), Linked Eurostat,¹³ and LinkedSpending,¹⁴ which contains government spendings from all over the world as Linked Data.

On enriching the RDF graphs coming out of the statistical tables, [25] annotates Web tables using class labels and relationships automatically extracted from the Web to augment the semantics and improve access. Integrated HTML tables are used to extend search aggregated results [2,27] and to insert Web table data into word processing software [2]. Enriching RDF graphs with missing temporal information has also been given attention in publishing historical data as RDF. For in-

stance, authors of [10] expose a knowledge graph that automatically integrates a spatio-temporal dimension from Wikipedia, GeoNames and WordNet data. Similarly, [21] proposes a generic approach for inserting temporal information to RDF data by gathering time intervals from the Web and knowledge bases. Authors in [8] focus on using the temporal aspect of Linked Data snapshots to keep track of the evolution of data over time. In our case, we extract temporal timestamps from legacy provenance information, as described in Section 3.6.

3. Data conversion and modeling

Our data conversion pipeline follows the diagram shown in Fig. 1. In the following sections we describe this pipeline in more detail.

3.1. Data conversion

In this section we describe the conversion process of the census tables from their original format to RDF.¹⁵ The dataset consists of 2,288 tables represented as spreadsheets in 507 Excel files. Each Excel file may contain one or several spreadsheets, but one spreadsheet always contains one single census table. An ex-

¹³See <http://eurostat.linked-statistics.org/>.

¹⁴See <http://linkedspending.aksr.org/>.

¹⁵All conversion source code is available at <https://github.com/CEDAR-project/Integrator/>.

RowHeader		HRowHeader	ColHeader	Data	Metadata	RowProperty
Gemeente	Nummer der beroepsklasse (NB: Romeinse cijfers)	Letter (Onderdeel beroepsklasse)	Regelnummer (NB: Arabische cijfers)	BENAMING van de onderdeelen der onderscheidene beroepsklassen, met de daartoe behoorende beroepen	Positie in het beroep (aangeduid met A, B, C of D)	
1	2	3				
Geboortejaren. leeftijd in j.	1878 en later. beneden 12 j.					
M	V	M	V			
O	O	O	O			
4	5	6	7			
Amsterdam	1			Aaedewerk, diamant, glas kalk, steenen, enz. Aardewerk en porcelein. Fabricage van aardewerk (incl. porselein, 1 terracotta, kachelbakkers, pottenbakkers enz.) 2 Fabricage van tabakspijpen Diamant, edelsteenen en fijne steensoorten. 3 Diamantslijpers (incl. verstellers) 4 Diamantslijpers (incl. verstellers) 5 Diamantslijpers (incl. verstellers) 6 Diamantslijpers (incl. verstellers) 7 Diamantsnijders	C A A B C D D	
17	3	128	11	5	12	

Fig. 2. One of the census tables of the dataset (occupation census of 1889, province of Noord-Holland). Colour markup is manually added and does not belong to the original data.

ample of such a table is shown in Fig. 2. A specific interpretation of the eccentric layout of these tables is necessary to generate RDF triples expressing exactly the same information. For instance, the bottom right figure in Fig. 2 should be read: there were 12 unmarried (*O* column) women (*V* column), 12 years old and born in 1878 (*12 1878* column) working as ordinary (row *D* in column *Positie in het beroep*, position in the occupation) diamond cutters (*Diamantsnijders* row) in the municipality of Amsterdam (column *Gemeente*, municipality). Consequently, this interpretation hampers a straightforward conversion of these tables, e.g. using well known generic community tools, to RDF (see Section 2). To this end, we developed TabLinker [16], a supervised Excel-to-RDF converter that relies on human markup on critical areas of these tables (see colors in Fig. 2). We define 6 markup styles that allow us to distinguish different cell roles (row headers, hierarchical row headers, column headers, data cells, metadata cells and row properties) within spreadsheets. With such markup, TabLinker can follow the same interpretation and generate meaningful RDF graphs across Excel files. The Integration pipeline shown in Fig. 1 uses the Integrator¹⁶ and TabLinker [16], which generates raw data according to our own table layout model instead of RDF Data Cube.

¹⁶See <https://github.com/CEDAR-project/Integrator/>.

3.2. Raw data

The Dutch historical censuses are multidimensional data covering a wide spectrum of statistics in population demography, labour force and housing situation. We choose RDF Data Cube (QB) as our goal data model to express the census data in RDF, since QB provides a means “to publish multi-dimensional data, such as statistics, on the web in such a way that they can be linked to related data sets and concepts” [5]. In QB, data points are called *observations*, primarily composed of a *measure* (e.g. “3 inhabitants”) and a set of *dimensions* qualifying that measure (e.g. “males”, “unemployed”, “in Amsterdam”). Dimensions can be arbitrarily combined to refer to unique observations in the cube.

However, the source tables lack critical information needed to generate a complete and sound QB dataset. Concretely, we miss mappings between dimensions with their corresponding values (e.g. it is said nowhere that column header *M* means *male* and relates to dimension *gender*, or that *O* means *unmarried* and relates to *marital status*). For this reason, we generate an agnostic RDF table layout representation as a first step, postponing the generation of proper RDF Data Cube.

After a 2 hour technical training, two people styled the 2,288 sheets of the dataset in 25 hours with the markup discussed in Section 3.1. Using such styles,

```

1 cedar:BRT_1889_08_T1-S0-K17 a tablink:DataCell ;
2   rdfs:label "K17";
3   tablink:dimension cedar:BRT_1889_08_T1-S0-A8 ;
4   tablink:dimension cedar:BRT_1889_08_T1-S0-K6 ;
5   tablink:dimension cedar:BRT_1889_08_T1-S0-J3 ;
6   tablink:dimension cedar:BRT_1889_08_T1-S0-K4 ;
7   tablink:dimension cedar:BRT_1889_08_T1-S0-K5 ;
8   tablink:dimension cedar:BRT_1889_08_T1-S0-B8 ;
9   tablink:dimension cedar:BRT_1889_08_T1-S0-C12 ;
10  tablink:dimension cedar:BRT_1889_08_T1-S0-E17 ;
11  tablink:dimension cedar:BRT_1889_08_T1-S0-F17 ;
12  tablink:value "12.0" ;
13  tablink:sheet cedar:BRT_1889_08_T1-S0 .

```

Listing 1. Raw RDF extracted for the cell K17 of the occupation census table of 1889, province of Noord-Holland.

TabLink first generates one `tablink:DataCell` for each data cell (i.e. cells marked as *Data* in Fig. 2), attaching its value (the actual population count) and the `tablink:sheet` the observation belongs to (a legacy table identifier, e.g. `BRT_1889_02_T1-S0`). Secondly, the observation is linked with all its corresponding column and row headers (i.e. cells marked as *RowHeader*, *HRowHeader*, and *ColHeader* in Fig. 2). An example is shown in Listing 1. Additionally, we create resources that describe the column and row headers, their types, labels, cell positions in the spreadsheets and hierarchical parent/child relationships with other headers (if any).

Because the result of this conversion stage is incomplete, due to the lack of further description of some dimensions and their mappings to standard values, codes and concept schemes, we call this the *raw* dataset conversion of the original Excel tables.

3.3. Integration rules as Open Annotations

To solve the missing dimension-value mappings shown in Listing 1, we annotate header cells using Open Annotation [22] with *harmonization rules* (see Listing 2). This is a manual process performed by experts. With such rules we can explicitly indicate the dimension to which a specific value belongs. Moreover, we can extend the description of such value (e.g. mapping “O” with “unmarried” and “V” with “female”) or map these values to dimensions that were not explicitly present in the original tables.

Some of these rules map the values extracted from the tables into standard *classification systems*. For instance, in order to query occupations consistently across the whole dataset, we map occupation dimension values (which are table dependent) to HISCO

```

1 cedar:BRT_1889_08_T1-S0-K4-mapping a oa:Annotation ;
2   oa:hasBody cedar:BRT_1889_08_T1-S0-K4-mapping-body ;
3   oa:hasTarget cedar:BRT_1889_08_T1-S0-K4 ;
4   oa:serializedAt "2014-09-24"^^xsd:date ;
5   oa:serializedBy
6     <https://github.com/CEDAR-project/Integrator> ;
7   prov:wasGeneratedBy
8     cedar:BRT_1889_08_T1-S0-mapping-activity .
9
10 cedar:BRT_1889_08_T1-S0-K4-mapping-body a rdfs:Resource ;
11   sdmx-dimension:sex sdmx-code:sex-F .

```

Listing 2. Mapping rules defined for *one* of the header cells associated to a data cell, in its corresponding annotation.

codes¹⁷ (Historical International Standard Classification of Occupations). We proceed similarly with other dimensions like historical religions, house types and historical municipalities in the Netherlands, using scripts and mappings done manually by experts (see Sections 4.1 and 4.2). We develop two tools to help experts on this process: LSD Dimensions, and TabCluster. LSD Dimensions [14] is an observatory of RDF Data Cube dimensions, codes, concept schemes and data structure definitions available now in the Linked Open Data cloud. It allows the reuse of these statistical resources by data owners and publishers. In case a specific concept scheme of interest is not available yet, we use TabCluster. TabCluster [15] is a concept scheme generator that leverages syntactic and semantic properties of non-standardized data cubes to assist data modelers on building concept schemes.

3.4. Harmonized RDF Data Cube

Using CONSTRUCT SPARQL queries, we process all the raw data produced by TabLink and apply all harmonization rules conveniently. As a result, we obtain refined, harmonized RDF Data Cube like shown in Listing 3. We generate a `qb:Observation` for each `tablink:DataCell`, and we link that observation to all its corresponding PROV triples.¹⁸

We also produce a `qb:DataStructureDefinition` (DSD) with all dimensions, attributes and measures used, and introduce several `qb:Slice` that group the observations by census type (VT, demography; BRT, occupations; and WT, housing) and year (from 1795 to

¹⁷See <http://historyofwork.iisg.nl/>.

¹⁸See “cube” module in <https://github.com/CEDAR-project/Integrator/>.

```

1 cedar:BRT_1889_02_T1-S0-K17-h a qb:Observation ;
2   maritalstatus:maritalStatus maritalstatus:single ;
3   cedar:occupationPosition cedar:job-D ;
4   cedar:population "12"^^xml:decimal ;
5   sdmx-dimension:sex sdmx-code:sex-F ;
6   prov:wasDerivedFrom cedar:BRT_1889_08_T1-S0-K17 ;
7   prov:wasGeneratedBy
8     cedar:BRT_1889_08_T1-S0-K17-activity .

```

Listing 3. Refined RDF Data Cube after applying harmonization rules in observation-attached OA annotations.

1971). The DSD can be browsed online,¹⁹ as well as the slices²⁰ and therefore all the observations.

3.5. Provenance

We implement provenance tracking with PROV [9] at all stages. We do this for a number of reasons. First, provenance allows us to ensure reproducibility of our conversion workflow. Second, it facilitates the debugging of all integration rules, since we can trace back all mappings, activities and entities involved in the generation of each qb:Observation. And third, we use it to meet the strong requirement of historians of being able to explain how every single harmonized value of the dataset is produced, back to the archived sources. For historians, ensuring independence and reliability of primary sources is fundamental, also in the Semantic Web [18].

For the TabLink generation of raw data cubes, we log a specific prov:Activity, recording task timestamps (prov:startedAtTime, prov:endedAtTime), its prov:Agent (prov:wasAssociatedWith) and the specific markup used via prov:used.

Similarly, during the execution of the mappings described as OA annotations we record an additional prov:Activity, making explicit the use of each specific mapping in the harmonization rules via prov:used.

3.6. Named graphs and URI policy

To organise the generated census triples we make them available in three different named graphs:²¹

- The *raw* data triples, as extracted from the original tables, are in <urn:graph:cedar:raw-data>.

¹⁹<http://lod.cedar-project.nl:8888/cedar/resource/harmonised-data-dsd>.

²⁰<http://lod.cedar-project.nl:8888/cedar/resource/harmonised-data-sliced-by-type-and-year>.

²¹Since we do not need them to be de-referenceable, we currently use URNs instead of URIs.

- All annotation mapping rules are contained in <urn:graph:cedar:rules>.
- The refined RDF Data Cube, produced after applying the mapping rules to the raw data, is located at <urn:graph:cedar:release>.

The resource URI naming policy is as follows: raw data cells are named following the schema cedar:(FILE-ID)-(SHEET-ID)-(CELL-ID), like cedar:BRT_1889_08_T1-S0-K17 (see Listing 1), where:

- (FILE-ID) is a legacy ID for the original Excel file, with the format (TYPE)-(YEAR)-(PART)-(VOLUME), e.g. BRT_1889_08_T1 refers to the occupation census (BRT) conducted in 1889, part 8, volume T1.
- (SHEET-ID) is an identifier of the sheet within a file, e.g. S0 for the first sheet, S1 for the second, etc.
- (CELL-ID) is an identifier of the cell within a sheet, e.g. K17 for the cell in column K, row 17.

The annotations containing the mapping rules associated to each header cell that affects a data cell follow exactly the same encoding, but adding the suffix “-mapping” to the resource. For example, cedar:BRT_1889_08_T1-S0-K4-mapping identifies the annotation containing the mapping rules for the header cell cedar:BRT_1889_08_T1-S0-K4.

Similarly, we identify the refined RDF Data Cube observations adding to the raw data URIs the suffix “-h”. For example, cedar:BRT_1889_08_T1-S0-K17-h identifies the qb:Observation we generate using the data cell cedar:BRT_1889_08_T1-S0-K17 as a basis and applying the mapping rules defined at the annotation cedar:BRT_1889_08_T1-S0-K17-mapping.

The three named graphs <urn:graph:cedar:raw-data>, <urn:graph:cedar:rules> and <urn:graph:cedar:release> contain a full conversion of the 2,288 census tables of the dataset. This collection is harmonized with a set of generic mapping rules, enabling the dataset to be queried under the RDF Data Cube schema. This collection is available at its own SPARQL endpoint.²² Additionally, we release a highly curated subset of this collection, called cedar-mini and contained in the named graphs <urn:graph:cedar-mini:raw-data>, <urn:graph:cedar-mini:rules> and <urn:graph:cedar-mini:release>. This subset is a high-quality curated harmonized version of the 59 most consulted and rel-

²²<http://lod.cedar-project.nl/cedar/sparql>.

Table 1

Datasets processed vs. datasets expected, with the total number of created RDF Data Cube observations

Description	Count
Number of datasets processed	1,358
Expected number of datasets	2,288
Total number of observations	6,800,175

evant tables of the collection. The cedar-mini subset is available at its own SPARQL endpoint.²³

4. Linked dataset description

In this section we describe the CEDAR dataset in more detail.²⁴ Table 1 shows some dataset statistics through its Data Structure Definition (DSD). Our conversion workflow is an ongoing process, since mapping rules in the observation annotations need to be manually curated. For this reason, we update these statistics every time we run the conversion workflow.²⁵ This allows us to keep track of what is left to map. Currently 6,800,175 observations are generated and linked to one qb:measureProperty (population), one qb:attributeProperty (unit of measure, number of persons), and nine qb:dimensionProperty: year of birth, sex, occupation position, belief, occupation, reference area, marital status, reference period, and census type.

Table 2 shows a summary of the different dimensions correctly mapped with standard codes into observations so far.

4.1. Internal links

The census tables often refer to variables and values with multiple synonyms: e.g. the value *female* for variable *sex* can be arbitrarily referred by *v*, *vrouw*, *vrouwen*, *vrouwelijk* or *vrouwelijk geslacht*.²⁶

In some variables this problem is straightforward to solve via the mappings we define as annotations, and we manually code mappings that cover all possible synonyms. This is the case for the variables **sex**, **marital status**, **occupation position** (i.e. rank class that a worker was assigned), **housing type situation**

Table 2

Dimensions of the dataset. The second column indicates how many observations in the dataset refer to such dimension. The third column indicates the proportion of observations referring to such dimension with respect to the total number of observations (6.8 M)

Dimension label	Occurrences	% obs.
belief	253,480	3.73%
censusType	4,642,360	68.27%
municipality	153,248	2.25%
maritalStatus	1,886,415	27.74%
occupation	328,790	4.84%
occupationPosition	8,120	0.12%
province	43,946	0.65%
refPeriod	4,642,360	68.27%
sex	3,801,431	55.90%

and **residence status**. Table 3 shows the correspondence between dimensions referenced in the observations, and the values (codes) they can get. The dimension *sex* is coded with `sdmx-dimension:sex`, and the codes `sdmx-code:sex-F` (female) and `sdmx-code:sex-M` (male) as values.²⁷ We mint our own URIs for dimensions *marital status* (`maritalstatus:maritalStatus`) and *occupation position* (`cedar:occupationPosition`). *Marital status* can get as value one of the codes `maritalstatus:single` (denoting single individuals), `maritalstatus:married` (married) or `maritalstatus:widow` (widows). Likewise, *Occupation position* can get as value one of the codes `cedar:job-D` (ordinary workers of the lowest rank, usually assigned to youth), `cedar:job-C` (ordinary workers with other lower-rank workers under their responsibility), `cedar:job-B` (foremen and other workers with many labour below their hierarchy) or `cedar:job-A` (directors or owners of businesses). The dimension *housing type situation* indicates the type of house inhabitants were counted in (occupied/empty houses, occupied/empty living ships, houses in construction), and *residence status* qualifies the status of the counted residents (present, legally registered and present, temporarily present, temporarily absent). The special predicate `cedar:isTotal` is used to mark observations that are aggregations over other observations. This distinction is important to avoid double counting when querying the dataset.

Other variables require a more complex schema of their possible values: for these QB suggests the use of concept schemes (also called *classification systems* in social history). The variable **house type**, which dis-

²³<http://lod.cedar-project.nl/cedar-mini/sparql>.

²⁴Unless stated, we refer to the collection with the full conversion of the 2,288 tables.

²⁵Full and regularly updated statistics can be found at <http://lod.cedar-project.nl/cedar/stats.html>.

²⁶*Vrouw* means *woman* in Dutch.

²⁷Some SDMX COG dimensions and codes are available in RDF at <http://purl.org/linked-data/sdmx/2009/dimension#> and <http://purl.org/linked-data/sdmx/2009/code#>.

Table 3

Dimensions linked from observations. The second column indicates whether we created (✓) or reused (×) the dimension. The third column indicates the possible dimension values in a concept scheme. The last column indicates how many observations contain such dimension value. Such references are expanded from a much smaller number of mapping rules, as shown in Table 5. Prefixes are described in Table 4

Dimension	New	Value/code in scheme	#Refs
cedar:houseType	✓	cedar:house-BewoondeHuizen	88,737
		cedar:house-BewoondeSchepen	28,573
		cedar:house-BewoondeWagens	4,221
		cedar:house-HuizenAanbouw	14,323
		cedar:house-OnbewoondeHuizen	51,599
		cedar:house-OverigeGebouwen	23,344
		cedar:isTotal	✓
cedar:population	✓	xsd:integer	710,462
cedar:residenceStatus	✓	resStatus:AltijdAanwezig	7,640
		resStatus:FeitelijkeAanwezig	81,625
		resStatus:JuridischAanwezig	220,293
		resStatus:TijdelijkAanwezig	119,373
		resStatus:TijdelijkAfwezig	55,733
		resStatus:WerkelijkTotaal	21,403
sdmx-dimension:refArea	×	From gg:10002 to gg:11447	692,491
sdmx-dimension:sex	×	sdmx-code:sex-M	220,661
		sdmx-code:sex-F	213,991

Table 4

Prefixes used in the dataset

Prefix	URI	Content
oa	http://www.w3.org/ns/openannotation/core/	Open Annotations vocabulary
prov	http://www.w3.org/ns/prov#	PROV ontology
dcat	http://www.w3.org/ns/dcat#	Data Catalog vocabulary
qb	http://purl.org/linked-data/cube#	RDF Data Cube (QB) vocabulary
sdmx-dimension	http://purl.org/linked-data/sdmx/2009/dimension#	QB dimensions
sdmx-code	http://purl.org/linked-data/sdmx/2009/code#	QB codes (dimension values)
sdmx-attribute	http://purl.org/linked-data/sdmx/2009/attribute#	QB attributes
cedar	http://lod.cedar-project.nl/vocab/cedar#	CEDAR terms: population, totals, variable descriptions
cedar-data	http://lod.cedar-project.nl:8888/cedar/resource/	CEDAR data points
resStatus	http://lod.cedar-project.nl/vocab/cedar-residenceStatus#	Residence status codes
maritalstatus	http://lod.cedar-project.nl/vocab/cedar-maritalstatus#	Marital status codes
gg	http://www.gemeentegeschiedenis.nl/amco/	Dutch historical municipalities (AMCO codes)
tablink	http://lod.cedar-project.nl/vocab/cedar-tablink#	TabLinker spreadsheet cell types
hisco	http://lod.cedar-project.nl/vocab/cedar-hisco#	Historical occupations (HISCO codes)

tinguishes military, civil, public and private buildings that were counted during the censuses, encodes building types in a taxonomic fashion. We manually build up this concept scheme²⁸ in a data-driven way, assisted by domain experts in social history.²⁹ We use the di-

mension cedar:houseType and an associated code list for this variable.

4.2. External links

Other variables, like **province**, **municipality**, **occupation** and **belief**, also need complex schemas or taxonomies to encode their values (see Section 4.1). We link to external datasets to standardize these variables.

²⁸ See concept scheme at <https://goo.gl/mt1dsn>.

²⁹ See [15] for an approach to build such taxonomies automatically.

Table 5

Type and number of mapping rules created per variable type. The second column links to the actual mapping files. The third column indicates how these mapping files were generated: either manually, by purely relying on expert knowledge (*expert-based*); or semi-automatically, with the aid of querying the raw data (*SPARQL*) or supported by *string similarity scripts*. The fourth column indicates the resulting number of mapping rules per file/variable. These mappings expanded into a much greater number of references to codes in concept schemes, as shown in Table 3

Variable	Mapping file	Generation	#Mapping rules
Age	https://goo.gl/5NIIqE	Expert-based; SPARQL	16,398
Belief	https://goo.gl/i1H2j4	Expert-based	582
City	https://goo.gl/poFcxo	Expert-based; string similarity script	42,294
Housing type	https://goo.gl/fdc0s8	Expert-based	3,484
Marital status	https://goo.gl/2rYLYu	Expert-based	10
Occupation	https://goo.gl/CUVSGc	Expert-based	21,851
Occupation position	https://goo.gl/y7NoYw	Expert-based	4
Province	https://goo.gl/yShX7w	Expert-based	18
Sex	https://goo.gl/ZtVS3z	Expert-based	10
Total	https://goo.gl/978YSy	Expert-based; SPARQL	38
Housing type situation	https://goo.gl/IEWfBf	Expert-based	22
Residence status	https://goo.gl/TRra0U	Expert-based	40

Province and *municipality* contain codes of Dutch provinces and municipalities from the past and are assigned as objects of predicates `sdmx-dimension:refArea`. Linking to GeoNames or DBPedia seems appropriate. However, Dutch provinces and municipalities suffered major changes during the historical censuses period. To address this, we issue links to `gemeentegeschiedenis.nl`.³⁰ `gemeentegeschiedenis.nl` is a portal that exposes standardized Dutch historical province and municipality names as Linked Open Data, based on the work done in the Amsterdamse Code (AC) [23]. 2,658,483 links are issued to provinces and municipalities in this dataset, based on previously existing manually curated mappings (see Table 5).

We follow a similar procedure to link values of the variable *occupation*. In this case, we rely on HISCO, which offers 1,675 standard codes for historical occupations. We issue 354,211 links to human-readable occupation description pages, also relying on existing manual mappings (see Table 5).

Other variables, like *belief* (religion), also need to be standardized by linking to standard classification systems. However, for these no proper historical classifications are available. In such cases, we create these classifications, either manually (relying on expert knowledge) or automatically (leveraging lexical and semantic properties [15]). In any case, we use mappings to these classifications to standardize the census values (see Table 5). We use such mappings to issue 256,952 links to historical religious denominations.

Table 3 shows a summary of the different dimensions mapped into observations so far, together with the codes associated to them, the number of times they are referenced, and whether they are created or reused from existing vocabularies. We also make available the RDF describing the created vocabularies,³¹ and we foresee a future reuse of these vocabularies by publishers of historical aggregated censuses of other countries. Table 4 lists all prefixes used. To standardize dimensions and their values, we create mapping rules and scripts; a summary of these is shown in Table 5 (string similarity scripts can be found online³²).

To achieve the fifth Linked Open Data star we produce links that connect the CEDAR dataset to other LOD datasets. Concretely, we issue links (see Fig. 3):

- To the Historical International Standard Classification of Occupations (HISCO) [24], which standardizes dimension values about historical occupations. The mappings were manually created by experts.³³
- To occupations in the ICONCLASS³⁴ system, offering alternative mappings to HISCO. The mappings were manually created by experts.³⁵

³¹See <https://github.com/CEDAR-project/Vocab>.

³²See <https://github.com/CEDAR-project/mapping-scripts>.

³³See mapping files at <https://github.com/CEDAR-project/DataDump/blob/master/mapping/Occupation.xls> and http://volkstellingen.nl/nl/onderzoek_literatuur/harmonisatie/beropen/index.html.

³⁴See <http://iconclass.org/>.

³⁵See mapping files at <https://raw.githubusercontent.com/CEDAR-project/DataDump/master/mapping/hisco-iconclass.csv>.

³⁰See <http://www.gemeentegeschiedenis.nl/>.

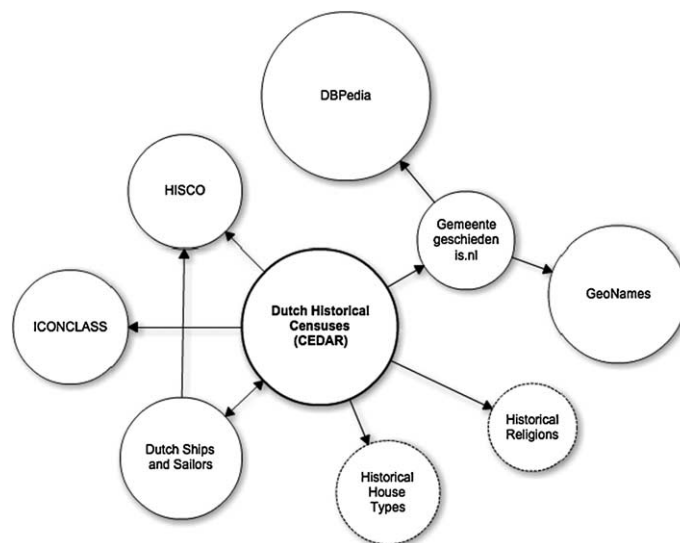


Fig. 3. Linked datasets to/from CEDAR.

- To URIs of gemeentegeshiedenis.nl (which point to resources in DBpedia and GeoNames) and the Amsterdamse Code [23]. These links standardize municipality names. These links were semi-automatically generated, by reusing existing mappings³⁶ into mapping scripts.³⁷
- To the Dutch Ships and Sailors dataset [6], linking interesting occupations and historical facts on Dutch maritime trade related to the census. These links were manually generated by an expert.

5. Impact and availability

5.1. Impact

Publishing the Dutch historical censuses as five-star Linked Open Data has a deep impact in the methodology that historians and social scientists have traditionally followed to study this dataset [1]. Due to the limitations of the old formats, the dataset could not be utilized to its full potential. To the date, most of the research based on the historical Dutch censuses focused on specific comparable years [3]. To utilize the full potential of the historical censuses researchers have iden-

tified harmonization of the data as a key aspect, which we implement as rules in `oa:Annotation` annotations. Previously, if researchers wanted to know e.g. the number of houses under construction in the Netherlands per municipality between 1859 and 1920,³⁸ they had to consult 47 different Excel tables and run into laborious data transformations. Moreover, keeping track of provenance of all performed operations was cumbersome and relied on data munging and delicate assumptions. By using explicit harmonization rules and links to standard classifications for occupations, municipalities, religions and house types, researchers can get answers to their queries in a blink of a time compared to the manual way of digging into disparate Excel tables. Table 6 shows the number of tables that users had to open and the number of cells they had to manipulate to answer a set of example queries. Most of these queries have been already manually investigated by social historians [3]. Hence, major milestones the dataset provides for History scholars are (a) speed-up of query answering; and (b) full provenance tracks of every data point down to the historical sources. Using the SPARQL endpoint, social scientists can retrieve data that gives support to hypotheses that previously could only be assumed. In addition, links to external datasets facilitate answering queries that users hardly could perform otherwise; for instance, links to gemeentegeshiedenis.nl and DBpedia allow to in-

³⁶See mapping files at <https://github.com/CEDAR-project/DataDump/blob/master/mapping/Cities.xls> and http://volkstelling.nl/nl/onderzoek_literatuur/harmonisatie/gemeenten/index.html.

³⁷See ‘CEDAR2gg’ at <https://github.com/CEDAR-project/mapping-scripts>.

³⁸Additional example queries at <https://github.com/CEDAR-project/Queries>.

Table 6

Example queries over the *cedar-mini* subset. For each query, we detail the number of tables that users had to open and the number of cells they had to manipulate in order to reach a query answer. Unless stated, reference periods cover from 1859 until 1920. SPARQL translations of these queries can be found at <https://github.com/CEDAR-project/Queries>

Query	#Tables	#Cells
Inhabited houses in Zuid-Scharwoude in 1889	1	1
Occupied houses and living ships per municipality	59	80,032
Legally registered and present inhabitants per municipality	34	23,086
Houses under construction	47	4,478
Empty houses	59	34,834
Temporarily present inhabitants in ships	35	4,255
Temporarily present inhabitants per municipality	47	74,462
Temporarily absent inhabitants per municipality	34	37,044
Temporarily present inhabitants in wagons	13	426
Number of houses according to their type, from 1859 until 1920	59	136,768
<i>Average</i>	38.8	39,538.6

stantly compare nowadays's population of Dutch municipalities with their historical figures, via SPARQL 1.1 federation. Moreover, dimension standardization enables new query solutions that were only possible through extensive manual work and expert knowledge.

As five-star Linked Open Data, the census dataset is open for longitudinal analysis, especially for a study of change. Being a major interest for historical research, the change in structures of classifications, meaning of variables and semantics of concepts over time, known as concept drift [26], is a fundamental topic to explore.

A set of tools built on top of the dataset is already available. For instance, social historians of the NLGIS project³⁹ query the endpoint to get historical census data and plot it in a map. Computational musicologists do research by linking the CEDAR dataset with their own historical singers database [11].

The dataset sums to other initiatives on publishing census data on the Web as RDF Data Cube.⁴⁰ To the best of our knowledge, ours is the first effort on publishing censuses with historical characteristics.

We have collected a number of SPARQL queries that we consider relevant for interested users.⁴¹ These are also available through the CEDAR dataset front-end.⁴²

The CEDAR dataset was used in the hackathon held during the 2014 CEDAR international symposium⁴³

with 11 attendees, and also in the 1st Digital History Datathon held at the VU University Amsterdam⁴⁴ with 13 attendees. The CEDAR dataset is listed as one of the datasets in the Challenge of the 2nd International Workshop on Semantic Statistics⁴⁵ (SemStats 2014), International Semantic Web Conference (ISWC 2014).

In addition, we log the usage of the dataset via any dereferenced URI or fired SPARQL query.

5.2. Availability

The CEDAR dataset, consisting of the raw Excel file conversions, the annotation mapping rules, and the harmonized RDF Data Cube, is served as Linked Open Data at <http://lod.cedar-project.nl/cedar/>. All URIs dereference via a Pubby installation on this server, which returns data formatted according to the requested format in the response header of HTTP requests. The dataset's SPARQL endpoint can be found at <http://lod.cedar-project.nl/cedar/sparql> (for the whole conversion) and <http://lod.cedar-project.nl/cedar-mini/sparql> (for the highly curated subset). All dataset dumps, including the original Excel files (with and without markup), mappings, and the converted RDF data can also be retrieved as bulk downloads at <https://github.com/CEDAR-project/DataDump>.

The creation and update of the dataset is done through a software package, the CEDAR Integrator,⁴⁶ developed for that purpose at the VU University Amsterdam and DANS under the LGPL v3.0 license.⁴⁷ The

³⁹See <http://www.nlgis.nl/>.

⁴⁰See cases for Italy, France, Australia and Ireland at http://www.istat.it/it/archivio/104317#variabili_censuarie, <http://goo.gl/R9iqQa>, <http://stat.abs.gov.au/> and <http://data.cso.ie/>.

⁴¹See <https://github.com/CEDAR-project/Queries>.

⁴²See <http://lod.cedar-project.nl/cedar/data.html>.

⁴³See <http://goo.gl/yfvUTI>.

⁴⁴See <http://cedar-project.nl/linkathon-at-the-vu/>.

⁴⁵See <http://semstats2014.wordpress.com/>.

⁴⁶See <https://github.com/CEDAR-project/Integrator>.

⁴⁷See <http://www.gnu.org/licenses/lgpl.html>.

Table 7

Breakdown of the coverage of dimension codes in the *cedar* and *cedar-mini* collections. The second column describes the dimension values being counted (“total” counts all raw dimension values and “unmapped” the ones that still need to be standardized). The third column links to an up-to-date SPARQL query performing the count. The last column indicates the number of dimension codes for each count

Collection	Count	Query	#Values
<i>cedar-mini</i>	Total codes	https://goo.gl/vX0fjM	188,958
<i>cedar-mini</i>	Total unique codes	https://goo.gl/m1TGO6	21,946
<i>cedar-mini</i>	Unmapped codes	https://goo.gl/l6m6QP	61,952 (32.79%), 67.21% mapped
<i>cedar-mini</i>	Unmapped unique codes	https://goo.gl/HymHX0	5,278 (24.05%), 75.95% mapped
<i>cedar</i>	Total codes	https://goo.gl/6mJehY	1,703,468
<i>cedar</i>	Total unique codes	https://goo.gl/WKp4oB	91,613
<i>cedar</i>	Unmapped codes	https://goo.gl/hu3v7h	1,150,950 (67.57%), 32.43% mapped
<i>cedar</i>	Unmapped unique codes	https://goo.gl/Pr1FNv	84,047 (91.74%), 8.26% mapped

dataset is regularly dumped to a GitHub repository.⁴⁸ Updates are performed in order to correct errors and incomplete mappings our experts detect when supervising statistical analyses⁴⁹ that we automatically generate during the conversion process (see Section 3.1). For long term preservation, the dataset is (and will continue being) deposited into DANS EASY,⁵⁰ a trusted digital archive for research data.

6. Discussion

In this paper we present the steps followed and the results achieved by CEDAR to transform a two-star (Excel conversions of scanned census tables) representation of the Dutch historical censuses into five-star Linked Open Data (harmonized census resources using URIs and linked to external concept schemes) as part of the Computational Humanities Programme⁵¹ of the Netherlands Royal Academy of Arts and Sciences.⁵²

We acknowledge a number of shortcomings in the dataset. Importantly, we are aware that the conversion is not complete. Although all observations reach the end of the pipeline, their mappings to standardized dimension values are incomplete. To address this, we follow two approaches: (a) we generate statistics during the conversion process;⁵³ and (b) we use SPARQL queries to analyse what mappings remain. Table 7 shows a comprehensive summary of the current outcome of such queries in both *cedar-mini* and *cedar*

collections. In *cedar-mini*, 75.95% of unique raw dimension values are currently mapped into standard codes. Consequently, 24.05% of values still need to be mapped. These values can be investigated also through SPARQL,⁵⁴ and mostly belong to very specific house types of the dimension *cedar:houseType*, some referring to unique historical buildings. In the *cedar* collection the missing mappings increase to 91.74%. This obviously includes the missing mappings of *cedar-mini*, but also:⁵⁵

- Occupational categories not included in the current occupation mapping files (dimension *cedar:occupation*), mostly referring to abstract categories in the occupations concept scheme.
- Values for the yet unmapped variable **age**. Age ranges are aggregated differently in each census edition, and mappings need to define additional interpolation rules in order to generate comparable data. Additionally, many values are redundant because of the existing duplicity between *age* and *year born* in the source data.
- Other geographical locations, like names of historical provinces and historical neighborhoods (dimension *sdmx-dimension:refArea*), for which no standard concept scheme exists.
- Historical religious denominations (dimension *cedar:belief*), for which no standard concept scheme exists.

With these approaches, we can quantify how far we are from completion and the work that still needs to be done on standardization. During the data generation we have issued temporal vocabularies and code

⁴⁸See <https://github.com/CEDAR-project/DataDump/> and <https://github.com/CEDAR-project/DataDump-mini-vt>.

⁴⁹See <http://lod.cedar-project.nl/cedar/stats.html>.

⁵⁰See <https://easy.dans.knaw.nl/>.

⁵¹See <http://www.ehumanities.nl/>.

⁵²See <http://www.knaw.nl/>.

⁵³See <http://lod.cedar-project.nl/cedar/stats.html>.

⁵⁴See <https://goo.gl/rQlkM2>.

⁵⁵The full list can be retrieved with the query <https://goo.gl/QUIOU5>.

lists for some variables that we will publish in separate data-hubs.⁵⁶ For instance, *belief* and *houseType* deserve their own Web spaces to allow other historical datasets to link to them. Linking the census observations to other datasets is another challenge.⁵⁷ Finally, the census tables contain a number of subtotals, totals and partial results at different levels of aggregation. We plan on checking the consistency of these aggregation levels automatically, spotting possible source errors.

Acknowledgements

The work on which this paper is based has been supported by the Computational Humanities Programme of the Royal Netherlands Academy of Arts and Sciences. For further information, see <http://ehumanities.nl>. This work has been supported by the Dutch national program COMMIT. We want to thank the reviewers for their thorough reviews, and Kathrin Dentler, Paul Groth and Andrea Scharnhorst for their valuable advice.

References

- [1] A. Ashkpour, A. Meroño-Peñuela and K. Mandemakers, The aggregate Dutch historical censuses: Harmonization and RDF, *Historical Methods: A Journal of Quantitative and Interdisciplinary History* **48**(4) (2015), 230–245. doi:10.1080/01615440.2015.1026009.
- [2] S. Balakrishnan, A. Halevy, B. Harb, H. Lee, J. Madhavan, A. Rostamizadeh, W. Shen, K. Wilder, F. Wu and C. Yu, Applying WebTables in practice, in: *Proc. of the Biennial Conference on Innovative Data Systems Research (CIDR 2015)*, M. Stonebraker, J. Gray and D. DeWitt, eds, 2015, <http://cidrdb.org/cidr2013/cidr2015proceedings.zip>.
- [3] O.W.A. Boonstra, P.K. Doorn, M.P.M. van Horik, J.G.S.J. van Maarseveen and J. Oudhof (eds), *Twee Eeuwen Nederland Geteld. Onderzoek met de digitale Volks-, Beroeps- en Woningtellingen 1795–2001*, DANS en CBS, The Hague, 2007, https://www.knaw.nl/nl/actueel/publicaties/twee-eeuwen-nederland-geteld/@download/pdf_file/Volkstelling_geheel_WEB_verkleind.pdf.
- [4] S. Capadisli, Towards linked statistical data analysis, in: *1st International Workshop on Semantic Statistics (SemStats 2013)*, ISWC, S. Capadisli, F. Cotton, R. Cyganiak, A. Haller, A. Hamilton, and R. Troncy, eds, Vol. 1549, CEUR, 2013, pp. 61–72, <http://ceur-ws.org/Vol-1549/article-06.pdf>.
- [5] R. Cyganiak, D. Reynolds and J. Tennison, The RDF Data Cube vocabulary, Technical report, W3C, 2014, <http://www.w3.org/TR/vocab-data-cube/>.
- [6] V. de Boer, J. Leinenga, M. van Rossum and R. Hoekstra, Dutch ships and sailors linked data cloud, in: *The Semantic Web – ISWC 2014, 13th International Semantic Web Conference*, Riva del Garda, Italy, October 19–23, 2014, P. Mika, T. Tudorache, A. Bernstein, C. Welty, C. Knoblock, D. Vrandečić, P. Groth, N. Noy, K. Janowicz and C. Goble, eds, LNCS, Vol. 8796, Springer, 2014, pp. 229–244. doi:10.1007/978-3-319-11964-9_15.
- [7] DERI, RDF Refine – A Google Refine extension for exporting RDF, Technical report, Digital Enterprise Research Institute, 2015, <http://refine.deri.ie/>.
- [8] V. Fionda and G. Grasso, Linking historical data on the web, in: *Proc. of the ISWC 2014 Posters and Demos Track, 13th International Semantic Web Conference (ISWC2014)*, M. Horridge, M. Rospocher and J. van Ossenbruggen, eds, Vol. 1272, CEUR-WS, 2014, http://ceur-ws.org/Vol-1272/paper_107.pdf.
- [9] P. Groth and L. Moreau, PROV-Overview, An overview of the PROV family of documents, Technical report, World Wide Web Consortium, 2013, <http://www.w3.org/TR/prov-overview/>.
- [10] J. Hoffart, F.M. Suchanek, K. Berberich and G. Weikum, YAGO2: A spatially and temporally enhanced knowledge base from Wikipedia, *Artificial Intelligence* **194**(28) (2013), 3161–3165. doi:10.1016/j.artint.2012.06.001.
- [11] B. Janssen, A. Meroño-Peñuela, A. Ashkpour and C. Guéret, Tracking down the habitat of folk songs, *eHumanities eMagazine* **4** (2015), <http://ehumanities.leasepress.com/emagazine-4/featured-article/tracking-down-the-habitat-of-folk-songs/>.
- [12] E. Kalampokis, A. Nikolov, P. Haase, R. Cyganiak, A. Stasiewicz, A. Karamanou, M. Zotou, D. Zeginis, E. Tambouris and K. Tarabanis, Exploiting linked data cubes with OpenCube toolkit, in: *Proc. of the ISWC 2014 Posters and Demos Track, 13th International Semantic Web Conference (ISWC2014)*, Riva del Garda, Italy, M. Horridge, M. Rospocher and J. van Ossenbruggen, eds, Vol. 1272, CEUR-WS, 2014, pp. 137–140, http://ceur-ws.org/Vol-1272/paper_109.pdf.
- [13] T. Lebo and J. McCusker, csv2rdf4lod, Technical report, Tetherless World, RPI, 2012, <https://github.com/timrdf/csv2rdf4lod-automation/wiki>.
- [14] A. Meroño-Peñuela, LSD dimensions: Use and reuse of linked statistical data, in: *Knowledge Engineering and Knowledge Management. EKAW 2014 Satellite Events, VISUAL, EKMI, and ARCOE-Logic*, Revised Selected, Linköping, Sweden, November 24–28, 2014, P. Lambrix, E. Hyvönen, E. Blomqvist, V. Presutti, G. Qi, U. Sattler, Y. Ding and C. Ghidini, eds, LNCS, Vol. 8982, Springer-Verlag, Berlin, Heidelberg, 2014, pp. 159–163. doi:10.1007/978-3-319-17966-7_22.
- [15] A. Meroño-Peñuela, A. Ashkpour and C. Guéret, From flat lists to taxonomies: Bottom-up concept scheme generation in linked statistical data, in: *Proc. of the 2nd International Workshop on Semantic Statistics (SemStats 2014)*, *International Semantic Web Conference (ISWC)*, S. Capadisli, F. Cotton, A. Haller, A. Hamilton, M. Scannapieco and R. Troncy, eds, Vol. 1550, CEUR,

⁵⁶Currently available vocabularies and code lists are available at <https://github.com/CEDAR-project/Vocab>.

⁵⁷See already issued links at <http://cedar-project.nl/linkathon-at-the-vu/>. Historical newspapers at <http://kranten.delfer.nl/> are other interesting data to link.

- 2014, pp. 72–77, <http://ceur-ws.org/Vol-1550/article-07.pdf>.
- [16] A. Meroño-Peñuela, A. Ashkpour, L. Rietveld, R. Hoekstra and S. Schlobach, Linked humanities data: The next frontier? A case-study in historical census data, in: *Proc. of the 2nd International Workshop on Linked Science (LISC2012), International Semantic Web Conference (ISWC)*, T. Kauppinen, L.C. Pouchard and C. Kessler, eds, Vol. 951, CEUR-WS, 2012, pp. 25–36, <http://ceur-ws.org/Vol-951/>.
- [17] A. Meroño-Peñuela, A. Ashkpour, M. van Erp, K. Mandemakers, L. Breure, A. Scharnhorst, S. Schlobach and F. van Harmelen, Semantic technologies for historical research: A survey, *Semantic Web – Interoperability, Usability, Applicability* 6(6) (2015), 539–564.
- [18] A. Meroño-Peñuela and R. Hoekstra, What is linked historical data? in: *Proc. of Knowledge Engineering and Knowledge Management. 19th International Conference, EKAW 2014*, Linköping, Sweden, November 24–28, 2014, K. Janowicz, S. Schlobach, P. Lambrix and E. Hyvönen, eds, LNAI, Vol. 8876, Springer-Verlag, Berlin, Heidelberg, 2014, pp. 282–287. doi:10.1007/978-3-319-13704-9_22.
- [19] T. Morris, T. Guidry and M. Magdinie, OpenRefine: A free, open source, powerful tool for working with messy data, Technical report, The OpenRefine Development Team, 2015, <http://openrefine.org/>.
- [20] E. Muñoz, A. Hogan and A. Mileo, DRETA: Extracting RDF from Wikitable, in: *Proc. of the ISWC 2013 Posters & Demonstrations Track, a Track Within the 12th International Semantic Web Conference (ISWC 2013)*, Sydney, Australia, October 23, 2013, E. Blomqvist and T. Groza, eds, Vol. 23, CEUR-WS, 2013, pp. 89–92, http://ceur-ws.org/Vol-1035/iswc2013_demo_23.pdf.
- [21] A. Rula, M. Palmonari, A.-C. Ngonga Ngomo, D. Gerber, J. Lehmann and L. Bühmann, Hybrid acquisition of temporal scopes for RDF data, in: *Proc. of the Semantic Web: Trends and Challenges, 11th International Conference, ESWC 2014*, Anissaras, Crete, Greece, May 25–29, 2014, V. Presutti, C. d’Amato, F. Gandon, M. d’Aquin, S. Staab and A. Tordai, eds, LNCS, Vol. 8465, Springer-Verlag, 2014, pp. 488–503. doi:10.1007/978-3-319-07443-6_33.
- [22] R. Sanderson, P. Ciccarese and H. Van de Sompel, Open Annotation Data Model, Technical report, W3C, 2013, <http://www.openannotation.org/spec/core/>.
- [23] A. van der Meer and O. Boonstra, *Repertorium van Nederlandse Gemeenten, 1812–2006, waaraan toegevoegd de Amsterdamse code*, DANS Data Guide 2, The Hague, 2006.
- [24] M. van Leeuwen, I. Maas and A. Miles, *HISCO: Historical International Standard Classification of Occupations*, Leuven University Press, 2002.
- [25] P. Venetis, A. Halevy, J. Madhavan, M. Paşca, W. Shen, F. Wu, G. Miao and C. Wu, Recovering semantics of tables on the web, *Proceedings of the VLDB Endowment (PVLDB)* 4(9), (June 2011), 528–538. doi:10.14778/2002938.2002939.
- [26] S. Wang, S. Schlobach and M.C.A. Klein, Concept drift and how to identify it, *Journal of Web Semantics: Science, Services and Agents on the World Wide Web* 9(3) (2011), 247–265. doi:10.1016/j.websem.2011.05.003.
- [27] M. Yakout, K. Ganjam, K. Chakrabarti and S. Chaudhuri, InfoGather: Entity augmentation and attribute discovery by holistic matching with web tables, in: *SIGMOD ’12 Proc. of the 2012 ACM SIGMOD International Conference on Management of Data*, C. Dyreson, F. Li and M.T. Özsu, eds, ACM, New York, NY, USA, 2012, pp. 97–108. doi:10.1145/2213836.2213848.