

Decomplexifying the network pipeline: a tool for RDF/Wikidata to network analysis

Julie M. Birkholz¹ and Albert Meroño-Peñuela²

¹WeChangEd, Department of Literary Studies, Ghent University,
Ghent, Belgium

²Department of Computer Science, Vrije Universiteit
Amsterdam, The Netherlands

Knowledge Graphs that use the Resource Description Framework language (RDF) as a knowledge representation paradigm are increasingly popular in Digital Humanities, and represent a valuable source of data for network analysis. However, digital scholars interested in network approaches over RDF graphs have to deal with complex workflows and frameworks in order to perform their analyses. These complexities exacerbate complications in reproducing and replicating their work. In this paper, we detail a proof of concept to combine popular libraries in RDF data management and network analysis in one single, publicly accessible Jupyter Notebook that enables a structured approach to network analyses of RDF graphs. What sets our work apart specifically is its flexibility in quickly re-running network analyses over slightly modified RDF graphs, and ensuring transparency in making the code visible. We explain this approach through two case studies: women editors in Europe in the 19th century, and provenance of the harmonization of the historical Dutch censuses (1795-1971). This approach affords the researcher to quickly, easily, efficiently and with increased reliability project and analyse networks from RDF.

1 Introduction

Linked Data is an increasingly common way to publish structured data in the Humanities ([de Boer et al., 2014](#), [Meroño-Peñuela et al., 2015](#), [Thornton et al., 2017](#)). As Tim Berners-Lee, the "creator" of the Semantic Web, described - Linked Data "provides a common framework that allows data to be shared and reused across application, enterprise, and community boundaries." ([\(W3C\), 2011](#)). Thus facilitating accessibility of knowledge on historical and cultural objects in a format readable by both humans and machines. For example, through standards such as the Resource Description Framework (RDF), natural language statements such as "George Orwell wrote 1984" can be expressed as a triple consisting of: a subject (:George_Orwell), a predicate (:wrote), and an object (:1984). This knowledge can be retrieved by machines

through a unique and global identifier (Uniform Resource Identifiers - URIs). This affords a networked archive, bringing together publicly available materials distributed in libraries, archives and museums; and thus allowing the researcher to integrate, and implement an unprecedented amount of often unstructured, siloed data, in lightning speed. Such an ontology or data model affords access, merging of information, and enrichment through efficiently linking of information on objects, entities and relations of collections to other collections.

Technically speaking data represented in the RDF language is structurally a graph. Thus it inherently allows us to infer relations, bundling any common affiliation between objects and attributes. From a research point of view this has led to a tendency to study RDF as a network. The study of networks and specifically the study of social networks has its roots in sociological theories where relationships form a part of the basis for understanding behavior (Durkheim, 1951, Simmel, 1955) where all actions are embedded in networks. (Granovetter, 1985) These relations – a set of edges, between nodes (entities) – define a network. These social networks reflect types of relations (e.g., a friendship tie in a friendship network or advice tie in an advice network). The study of networks, and in particular social networks, have been and are on the rise, providing explanations for relational and systematic phenomena (Borgatti and Foster, 2003), as it moves beyond explanations based on individual factors. For example not that someone's age explains their success, but rather the structure of their social network. (Granovetter, 1985)

The identification of networks is often thought of as a laborious task. It is traditionally done in many fields by searching through archival sources to identify nodes and edges, and reshaping data that is often not collected as relational, but from which one can infer relations. This entails integrating, and implementing a large amount of often unstructured, siloed and incomplete data to reconstruct relations between nodes and edges. Thus information about relations where a social network can be inferred from RDF provides a great advantage for exploring social networks embedded in this data. Networks can efficiently be reconstructed with the development of specific SPARQL queries to reflect different lenses of relations. For example, generating networks of different time periods, of different types of relations, with different boundaries (looking at relations of one city versus one country, or a neighborhood to a street) over the same data source.

Modeling data as networks affords the implementation of network analysis. Network analysis –the method used to analyze relations– provides a lens to investigate these diverse complex relational dynamics to examine structure, content or function. For social networks, which are the focus of the examples we provide in this paper, the structure of networks and positions of actors in these structures are seen as proxies for understanding social structure (Burt, 1980, Coleman, 1988).

The analysis of networks from RDF is largely done with a pipeline of tools (i.e. (Gil and Groth, 2011, Groth and Gil, 2011)). This starts with a data source, and the tools necessary for querying the specific data and relations. For example one workflow may be: the Wikidata Query Service, which allows one to query linked data in the Wikisphere through a SPARQL query, and can be exported in a number of formats; or the data might be stored in a database and is extract-able as a JSON(-LD) file. Moving from these file types, this relational data needs to be converted into a file type that is readable by a network analysis software. Typical network analysis software use a range of inputs depending on the program. The two most commonly used user friendly network analysis and visualization programs with a graphic user interface

are Gephi¹ and UCIne². These programs allow the implementation of various types of input files; for example: .csvs, matrices and DL files, as well as program specific files. Then pending the required analysis there are a number of export options to further reuse these results as data. This, for example, could include analysing network measures and considering them as a variable in a statistical model in a program such as SPSS, or R. Thus the current workflow approaches for working with network data from RDF requires researchers to work through multiple programs to specify queries, extract networks and export data as matrices, and implement network analysis tools to investigate graphs.

In addition, in building such a pipeline we lose sight of the hermeneutics of the research objects. (Gibbs and Owens, 2013) Researchers are often faced with black boxed tools that limit their understanding of the projection, generation, analysis or reformatting that occurs with each step. With each use of an additional program, algorithm or command, the data gets re-"massaged" and shaped. This further becomes an issue, when the development of such a pipeline is a technical adversary for domain experts (e.g. historians, literary scholars) with (traditionally) limited technical knowledge; but also for researchers with specific expertise in RDF or networks. Thus, we argue there is a need, within the DH community, to reduce this RDF-to-network analysis pipeline without creating another domain or research question specific tool, and while maintaining oversight over the process from RDF-to graph-to network analysis.

To address these issues, we propose the use of a Jupyter notebook that integrates the Python packages: RDFLib³ with NetworkX^{4,5}. This results in a reusable workflow that allows network analyses over RDF data to be more accessible, flexible, transparent and iterative. This is due to that increases the reliability in exploring all the possible social networks within the available RDF, as well as increases the speed, ease, and efficiency of the necessary steps of RDF to network analysis. What specifically sets our work apart from previous workflows is its flexibility in quickly re-running network analyses over slightly modified RDF graphs, while maintaining the code visible for transparency and learning. We outline this pipeline through two case studies:

1. a case of the social networks of 19th century women editors in Europe available on Wikidata, and
2. provenance of the harmonization of the historical Dutch censuses (1795-1971)

to explain how it can be useful for humanities research.

2 Method

The notebook consists of five "cells", which are actionable code blocks, shown here in Figure 1. The output of all these processes can be selected and copy-pasted for further reuse in graph processing frameworks or directly in reports or papers.

1 <https://gephi.org/>
2 <https://sites.google.com/site/ucinetsoftware/home>
3 <https://github.com/RDFLib/RDFLib>
4 <https://NetworkX.github.io/>
5 The full notebook is available at <https://github.com/descepolo/rdf-network-analysis/blob/master/rdf-network-analysis.ipynb>. A Google Colaboratory version of the notebook is also available, which makes it executable on the web with no need of local installation: <https://colab.research.google.com/github/descepolo/rdf-network-analysis/blob/master/rdf-network-analysis.ipynb>

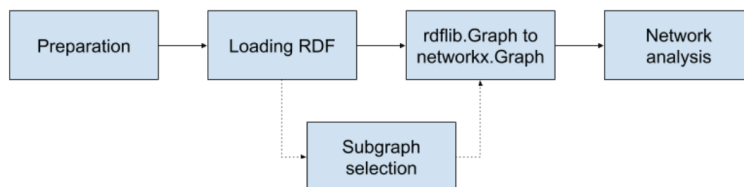


Figure 1: Workflow of the RDF Network Analysis Jupyter notebook

2.1 Preparation

As a first step the notebook loads the relevant packages - RDFLib and NetworkX. RDFLib is a Python package for working with RDF that includes parsers and serializers for RDF/XML, N3, NTriples, N-Quads, Turtle, TriX, RDFa and Microdata; a graph interface; store implementations for in memory storage and persistent storage on top of the Berkeley DB; and a SPARQL 1.1 implementation. (Krech, 2006) This facilitates a flexible environment for loading and manipulating RDF graphs. Then the user is prompted to input the full path to an RDF graph to load the RDF graphs. This can be any local or online RDF file.

2.2 Subgraph Selection

Users select a specific network in the RDF graph. The efficient aggregation of different snapshots of the networks can be achieved through a SPARQL query. SPARQL is a Semantic Web query language for databases which enable the ability to retrieve and manipulate data RDF specifically (Segaran et al., 2009).

2.3 From RDFLib to NetworkX

In order to generate a network, this RDF needs to be converted into a matrix. This is accomplished through a conversion of RDFLib.Graph to NetworkX.Graph. This prepares a file of the identified graph for analysis in NetworkX.

The Python library NetworkX enables the analysis of networks of around 10 million nodes and 100 million edges. (Hagberg and Conway, 2010) It is ideal for use for digital humanities as it affords the use of many types of networks, including directed graphs, and graphs with and self loops; while not maintaining strict object functions. (Hagberg et al., 2008) This implies that in the case of RDF which may have many and multiple types of networks embedded in the triples it will model anything that is structured as a matrices. This could include networks that we do not discuss here in this paper such as affiliation or two-mode networks, semantic networks and so forth. Thus the tool, which operates in the more general space of RDF models, does not limit the boundaries of inspection by imposing specific network models, leaving this choice to the user.

2.4 Network Analysis

Networks can be represented as graphs where positions and structures are systematically analyzed. (Wasserman et al., 1994) These principles originate from graph theory, which provides mathematical descriptions of characteristics. (Van Steen, 2010)

The networks can then be analyzed in NetworkX considering a number of characteristics of the network, as well as statistical analyses, see Table 1 Proposed Network

Characteristics. We have selected a standard, non-exhaustive, set of one-mode complete network measures. This is to establish the proof of concept, of course in practice any network measure that is included in NetworkX could be implemented in this notebook, for example measures of community detection, to other measures of centrality. It is not the goal of this paper to explain the operationalization of theoretical concepts to network measures, but it should be considered by humanities researchers in deciding on applicable measures to include in their research. For a more exhaustive list and explanation of network measures see (Wasserman et al., 1994).

Following this selection the network analysis is run and the results are printed, as well as a basic visualization which serves for the researcher to confirm a first accuracy check of the network, e.g. were the correct node and edges selected?; does something look strange or potentially missed in the query?, that can now be amended.

Network Concepts	Network measures
network size	total number of nodes, and the average number of edges
power centrality	nodal position: e.g. degree centrality, betweenness, and eigenvector centrality (Freeman, 1978)
density	a value of the proportion of all possible ties that are present

Table 1: Network Characteristics.

3 Case Studies

In this Section we validate our approach using two different case studies for the Digital Humanities: the social networks of women editors in Europe in the 19th century; and the provenance graphs of harmonization transformations performed in the Dutch historical censuses. The use of these cases are to demonstrate the use of the notebook, not a network study with elaborated research questions and operationalized network measures.

3.1 Women Editors in Europe in the 19th century

The 19th century in Europe, was one of the onset and rise of industrialization, altered the socioeconomic and cultural norms influencing the movement of people through advancements in train infrastructure and technologies in food and consumer goods, and investments in education throughout Europe. This also led to an increasing advancement of women's rights and positions in society. The ERC "Agents of Change: Women Editors and Socio-Cultural Transformation in Europe, 1710-1920" (acronym WeChangEd) directed by Marianne Van Remoortel and based at the Department of Literary Studies, Ghent University, Belgium (project Agents of Change: Women Editors and Socio-Cultural Transformation in Europe 2015), questioned how the press and periodical editorship in particular enabled women to take a prominent role in public life, to influence public opinion and to shape transnational processes of change. To facilitate the collection of biographical records, and archival evidence of women editors in Europe a Linked Data model was developed (Schelstraete and Van Remoortel, 2019).

This model afforded the cataloguing and tracing of different social networks in which the women participated.

This resulted in a large and growing database which includes 1700+ persons, 1600+ periodicals and 200+ organizations, as well as biographical information of these entities and relations between them, as identified through archival research. This data is available as the WCD Database, as subsets of data stored as .csv (Van Remoortel 2020). In April 2020, the WCD database was imported to Wikidata (Thornton et al. 2020) to facilitate the reuse and integration of this information with other Linked Open Data sources. The WeChangEd data can be identified in Wikidata through the unique property instance of WeChangEd ID P7947, see - <https://www.wikidata.org/wiki/Property:P7947>. This resulted in 3661 instances of data which comprises people, periodicals, organizations, as well as records of the relationships between these three entities, biographical information about these instances, and so forth. The complete dataset can be found via a Wikidata Query Service via - <https://w.wiki/QiQ>.

Identifying historical social networks is a laborious task, thus having the information on relations in Wikidata, and specifically as RDF, allows the researcher to explore historical social networks of the past in a more valid and flexible manner. The validity is increased, as the information is shared with the community, where it can be cross-checked, questioned, and enriched through the edit functions of Wikidata. As we show here through this example, the flexibility is affording by this pipeline.

In exploring how a researcher can identify social networks of these editors we display here three examples of projecting personal relationships of female editors between individuals as identified within the WeChangEd dataset. To identify these relationships we developed three SPARQL queries for the Wikidata Query Service, which we detail here below, and are also available at: <https://w.wiki/Qtr>, <https://w.wiki/QiQ>, and <https://w.wiki/QcS>, respectively. Using these graphs as input for the method described in Section 2, we convert these graphs to a NetworkX file and the network analysis is executed. We implement this query in the notebook resulting in three different network projections, and reflect on the implications for digital humanities researchers in compiling social networks from the past.

The first network represents a query on the entire WCD dataset, to identify kinship relations, this includes any identified siblings, parents, unmarried partner, spouse, or children of female editors <https://w.wiki/Qtr> (see Quelltext 1). This results in a network of all female editors and their relationships as identified in Wikidata, where nodes are individuals and edges or ties of represent a personal relationship, see Figure 2.

```

1 SELECT DISTINCT ?item ?o ?itemLabel ?sibling ?spouse ?partner ?father
  ?mother ?child
2 WHERE
3
4 {
5 # find occupation editors
6 ?item wdt:P106 wd:Q1607826.
7 ?item wdt:P7947 ?o .
8
9 # that are female
10 ?item wdt:P21 wd:Q6581072.
11
12 # that have a birth and death date
13 ?item wdt:P569 ?birthDate.
14 ?item wdt:P570 ?deathDate .
15
16 # with kinship: sibling
17 OPTIONAL { ?item wdt:P3373 ?sibling .}
18 # with kinship: spouse
19 OPTIONAL { ?item wdt:P26 ?spouse .}
20 # with kinship: unmarried partner
21 OPTIONAL { ?item wdt:P451 ?partner .}
22
23 # with kinship: father
24 OPTIONAL { ?item wdt:P22 ?father .}
25 # with kinship: mother
26 OPTIONAL { ?item wdt:P25 ?mother .}
27 # with kinship: child
28 OPTIONAL { ?item wdt:P40 ?child .}
29
30 # labels
31 SERVICE wikibase:label { bd:serviceParam wikibase:language
32 "[AUTO_LANGUAGE],en". }
33
34 } ORDER BY ?birthDate ?deathDate

```

Quelltext 1: SPARQL query for all female authors and their kinship relations

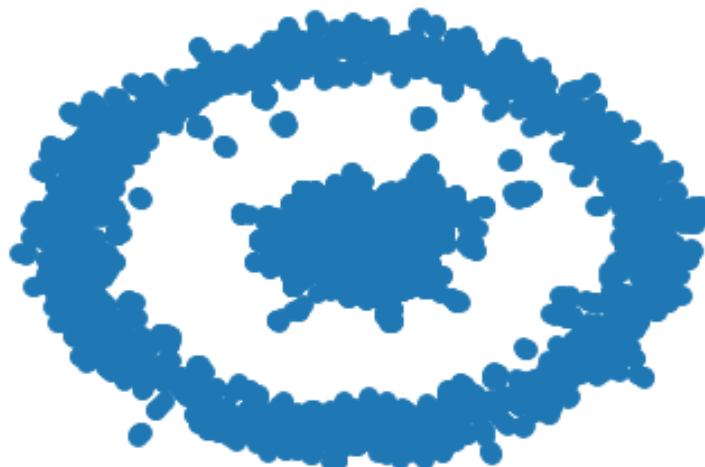


Figure 2: Network of editors

In this second selection we aim to show, how to refine the query, to select a more bounded set of nodes. This is a bounded selection of relations from within the WCD dataset but specifically of 19th century British female editors and their kinship relations, this includes any identified siblings, parents, unmarried partner, spouse, or children:

```

1 SELECT DISTINCT ?item ?o ?itemLabel ?sibling ?spouse ?partner ?father
  ?mother ?child
2 WHERE
3
4 {
5 # find occupation editors
6 ?item wdt:P106 wd:Q1607826.
7 ?item wdt:P7947 ?o .
8
9 # that are female
10 ?item wdt:P21 wd:Q6581072.
11
12 # that have a birth and death date
13 ?item wdt:P569 ?birthDate.
14 ?item wdt:P570 ?deathDate.
15
16 # that is British
17 ?item wdt:P27 wd:Q174193.
18
19 # with kinship: sibling
20 OPTIONAL { ?item wdt:P3373 ?sibling .}
21 # with kinship: spouse
22 OPTIONAL { ?item wdt:P26 ?spouse .}
23 # with kinship: unmarried partner
24 OPTIONAL { ?item wdt:P451 ?partner .}
25
26 # with kinship: father
27 OPTIONAL { ?item wdt:P22 ?father .}
28 # with kinship: mother
29 OPTIONAL { ?item wdt:P25 ?mother .}
30 # with kinship: child
31 OPTIONAL { ?item wdt:P40 ?child .}
32
33 # only active in the 19th century
34 FILTER ( ?birthDate >= "1800-01-01T00:00:00Z"^^xsd:dateTime &&
35 ?deathDate <= "1898-12-31T00:00:00Z"^^xsd:dateTime )
36
37 # labels
38 SERVICE wikibase:label { bd:serviceParam wikibase:language
39 "[AUTO_LANGUAGE],en". }
40
41 } ORDER BY ?birthDate ?deathDate

```

Quelltext 2: SPARQL query for relationships of British female editors of periodicals in the 19th century in Wikidata

<https://w.wiki/QnA> (see Quelltext 2). This results in a network of the personal relations of 19th century British female editors, where nodes are individuals and edges are relationships, see Figure 3. This network is a subset of the larger graph, but with parameters of time - editors living during the 19th century, and place - what was then the United Kingdom of Great Britain and Ireland.

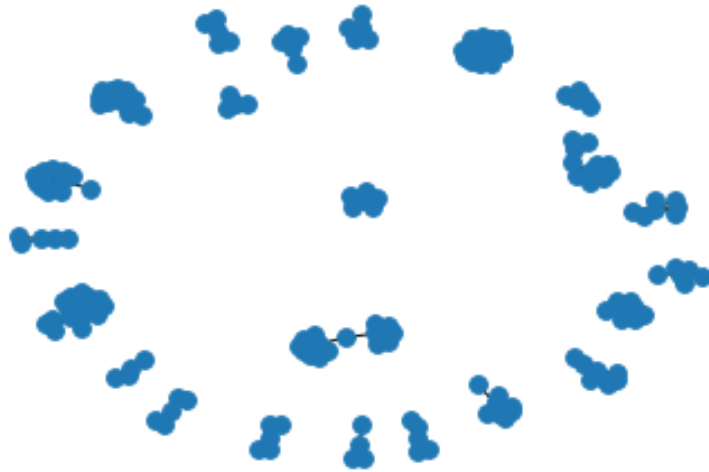


Figure 3: Network of 19th Century British female editors

The third case, aims to represent a different subset of the data, that is looking at relationships between editors based on language, instead of a geographical or political space. This selection represents a query of 19th century German speaking female editors and their kinship relations, this includes any identified siblings, parents, unmarried partner, spouse, or children: -<https://w.wiki/QnB> (see Quelltext [3](#)).

```

1 SELECT DISTINCT ?item ?o ?itemLabel ?sibling ?spouse ?partner ?father
2 ?mother ?child
3 WHERE
4 {
5 # find occupation editors
6 ?item wdt:P106 wd:Q1607826.
7 ?item wdt:P7947 ?o .
8
9 # that are female
10 ?item wdt:P21 wd:Q6581072.
11
12 # that have a birth and death date
13 ?item wdt:P569 ?birthDate.
14 ?item wdt:P570 ?deathDate.
15
16 # that speaks German
17 ?item wdt:P1412 wd:Q188.
18
19 # with kinship: sibling
20 OPTIONAL { ?item wdt:P3373 ?sibling .}
21 # with kinship: spouse
22 OPTIONAL { ?item wdt:P26 ?spouse .}
23
24 # with kinship: father
25 OPTIONAL { ?item wdt:P22 ?father .}
26 # with kinship: mother
27 OPTIONAL { ?item wdt:P25 ?mother .}
28 # with kinship: child
29 OPTIONAL { ?item wdt:P40 ?child .}
30
31 # only active in the 19th century
32 FILTER ( ?birthDate >= "1800-01-01T00:00:00Z"^^xsd:dateTime &&
33 ?deathDate <= "1898-12-31T00:00:00Z"^^xsd:dateTime )
34
35 # labels
36 SERVICE wikibase:label { bd:serviceParam wikibase:language
37 "[AUTO_LANGUAGE],en". }
38
39 } ORDER BY ?birthDate ?deathDate
40

```

Quelltext 3: SPARQL query for relationships of German speaking editors of periodicals in the 19th century in Wikidata

This results in a network of relations of German speaking female editors as identified on Wikidata. Selecting German speaking instead of a specific empires and or nation-state provides a broader query for identifying possible interactions between the German-speaking community in the 19th century. This results in a network of individuals as nodes and edges as relations, see Figure 4

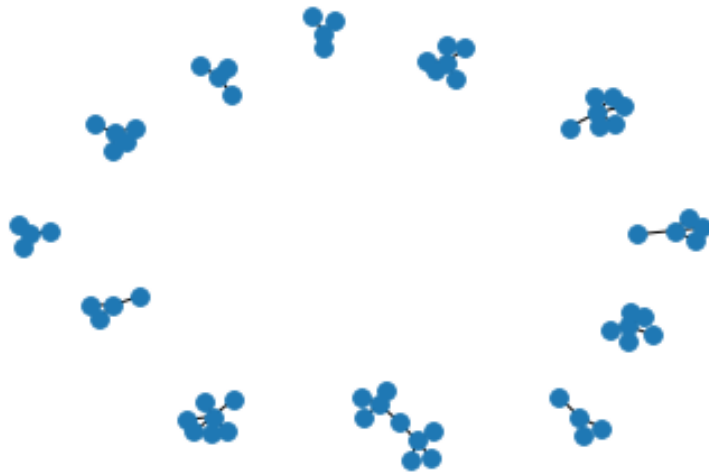


Figure 4: Network of 19th Century German speaking female editors

The complete results for two specific social networks of 19th century British female editors and 19th century German speaking female editors can be found in detail in the appendix. The results show the network analysis on connected components or groups of connected individuals, most central nodes, and communities. A researcher can use these results, combined with the visualizations to further explore these relations either returning to archival materials to investigate previously understudied relations, or further analyse the structure and positions within these network to explain social capital of the periodicals the editors edited or kinship relations.

These three examples from within the WCD dataset on Wikidata display the flexibility of this approach in moving through a dataset, to generate social networks. This notebook, in contrast to other workflows allows researchers to consider aspects of space, time, place and other parameters of the data within a few steps and seconds; where the researcher can move and back and forth between the raw data, the query, the network projection, and the analysis, to compile the most suitable, reliable graph from the available data. It serves as both an efficient approach to explore the social relations within a dataset, as well as to validly and reliably generate a social network and conduct social network analysis of the networks from RDF.

3.2 CEDAR: Harmonization Provenance of the Dutch Historical Censuses (1795-1971)

The Dutch historical censuses were collected in the Netherlands in the period 1795–1971, in 17 different editions, once every 10 years. The government counted all the country’s population, door-to-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, this exhaustive collection stopped due to social opposition, and the government switched to municipal registers and sampling (Ashkpour et al., 2015). Various projects have digitized the resulting census data (CBS; IISH; Data Archiving and Networked Services⁶; DANS; and

⁶ See <http://www.dans.knaw.nl/>

the Netherlands Interdisciplinary Demographic Institute⁷ NIDI), and have manually translated them into a collection of 507 Excel spreadsheets and 2,288 census tables.⁸ The CEDAR project⁹ takes these spreadsheets as input, and produces a Knowledge Graph of 6.8 million statistical observations (Meroño-Peñuela et al. 2015) many of which went through an harmonization process to satisfy the standardization needs of historians for their querying (Ashkpour et al. 2015).

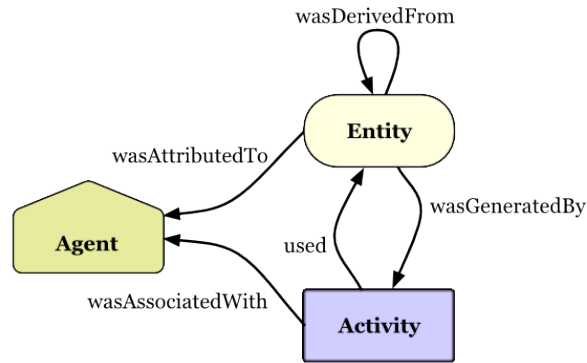


Figure 5: Provenance model of W3C PROV (Lebo et al. 2013).

In this case study, we use the CEDAR Knowledge Graph (Meroño-Peñuela et al. 2015) with our proposed approach to explain how to a researcher can consider network similarities and differences between various historical census data points and their *provenance* information. Historians are particularly interested in the transformation and manipulations that occurred in this harmonization process in generating these data points; as this signals their correctness and hence its reliability. Fortunately, the CEDAR Knowledge Graph documents the harmonization transformations of all data points using the W3C PROV standard (Lebo et al. 2013). This standard models provenance as the interactions between various *entities* (the objects subject to transformations, i.e. the census data points), *activities* (the transformation processes themselves, i.e. the harmonization rules) and *agents* (the persons or programs commanding the transformations) as shown in Figure 5

We select two arbitrary observations of the census, VT_1859_01_H1-S8-J647-h (observation 1, o_1) and VT_1920_01_T-S0-R10108-h (observation 2, o_2), and their corresponding provenance traces with the query shown in Listing 4 against the CEDAR SPARQL endpoint¹⁰ We use the graphs returned by this query as input for the notebook¹¹

We use the provenance graphs of o_1 and o_2 as input for the method described in Section 2 We execute the preparation block; we use the query of Listing 4 as subgraph selection; we execute the network conversion block; and finally we execute the network analysis block. The output networks as plotted by the notebook are shown in Figure 13 We can observe that for both cases the network is 2-star shaped, with the nodes repre-

7 <http://www.nidi.knaw.nl/en/>
 8 <http://volkstellingen.nl/>
 9 <https://www.cedar-project.nl/>
 10 <https://api.druid.datalegend.net/datasets/datalegend/CEDAR-S/services/CEDAR-S/sparql>

11 The input graphs are also available at <https://github.com/albertmeronyo/rdf-network-analysis/blob/master/uc1.nt> and <https://github.com/albertmeronyo/rdf-network-analysis/blob/master/uc2.nt>

```

1 CONSTRUCT {
2   ?obs ?obs_p ?obs_o .
3   ?act ?act_p ?act_o .
4 } WHERE {
5   VALUES ?obs { :VT_1859_01_H1-S8-J647-h :VT_1920_01_T-S0-R10108-h}
6   ?obs prov:wasGeneratedBy ?act .
7   ?obs ?obs_p ?obs_o .
8   ?act ?act_p ?act_o .
9 }

```

Quelltext 4: SPARQL query for the harmonization provenance graphs of two census observations.

senting the observation and the activity at the center of these stars and various nodes describing their properties, as expected. One edge (`prov:wasGeneratedBy`) connects these two nodes. An noticeable difference is that while o_1 (Fig. 11) is transformed by 6 different harmonization rules, o_2 (Fig. 12) is only affected by 3. This can provide interesting insights for historians, who may be keen to examine statistical observations that have been subject to a higher number of transformations (and therefore more prone to errors) and the relations of these transformations to their immediate context. In this sense, visualizing these network contexts can be a powerful tool for interpretation.

Additionally, Figure 10 shows the histograms for degree and eigenvector centrality drawn by the notebook for both graphs. This is a more aggregated view on the networks, showing similar behaviour for o_1 and o_2 (due to the structural similarity of provenance graphs) but also interesting differences. For example, o_1 eigenvector centrality shows a more normal distribution due to the higher variety of node influence in a more varied network. The remaining network statistics can be found upon the execution of these two examples in the notebook at <https://github.com/descepolo/rdf-network-analysis/blob/master/rdf-network-analysis.ipynb>.

4 Conclusion and Future Work

In this paper we have detailed how we have proposed to combine popular libraries in RDF data management and network analysis in one single, publicly accessible Jupyter Notebook that enables a structured approach to network analyses of RDF graphs. With the proposed Jupyter notebook we have developed a transparent and iterative tool for RDF to network in research. The open code and user-friendliness of the notebook ensures flexibility for users in implementing different aspects of the two libraries that we did address here in this demonstration. In addition, we have demonstrated through the tool and use cases how this affords the reuse and accessibility for non-technical scholars of RDF, as well as increase the efficiency and flexibility of use for generating networks from RDF. This approach facilitates the study of diverse types of networks from RDF and thus study of relational phenomenon in the Humanities and beyond.

In addition, this approach proves, contrary to the trend in the digital humanities, that we do not need a new network software that converts diverse file types to make fundamental improvements on both the quality of the networks used in research, as well as the the analysis of networks. Rather, as we presented, a fundamental rethinking of how data on social networks is structured, manipulated and pushed through a pipeline is needed to efficiently generate, project and evaluate networks. This approach increases the flexibility, compared to traditional network workflows- where the analyst would prepare a matrix for each projection of a network, go back to source material every time to reshape the data and networks based on different periods, or parameters

(e.g. variables such as country of birth, gender, language of entities), and push it through the workflow. Such an approach reduces the technical adversary of knowledge on RDF and network analysis, while avoiding a black boxed software, as well as retains a hermeneutic approach to the source data, allowing the researcher to iteratively and efficiently requery, reshape and reanalyze the networks embedded in RDF.

Acknowledgements

The authors would like to thank Prof. dr. Paul Groth for his feedback on this work but also for his ongoing support of all of our crazy ideas. The first author would like to thank her partner- Radboud Reijn for being so patient, and secondly the continued support from Prof. dr. Marianne Van Remoortel in proposing and believing in a DH project for periodical studies. The work of the first author was supported by the H2020 European Research Council under the ERC Starting Grant agreement no. 639668. This work has been partly supported by the Dutch national project CLARIAH.

References

- Ashkpour, A., A. Meroño-Peñuela, and K. Mandemakers
2015. The aggregate Dutch historical censuses: Harmonization and RDF. *Historical Methods: A Journal of Quantitative and Interdisciplinary History*, 48(4):230–245.
- Borgatti, S. P. and P. C. Foster
2003. The network paradigm in organizational research: A review and typology. *Journal of management*, 29(6):991–1013.
- Burt, R. S.
1980. Models of network structure. *Annual review of sociology*, 6(1):79–141.
- Coleman, J. S.
1988. Social capital in the creation of human capital. *American journal of sociology*, 94:S95–S120.
- de Boer, V., M. van Rossum, J. Leinenga, and R. Hoekstra
2014. Dutch ships and sailors linked data. In *International Semantic Web Conference*, Pp. 229–244. Springer.
- Durkheim, E.
1951. Suicide: A study in sociology (ja spaulding & g. simpson, trans.). *Glencoe, IL: Free Press. (Original work published 1897)*.
- Freeman, L. C.
1978. Centrality in social networks conceptual clarification. *Social networks*, 1(3):215–239.
- Gibbs, F. and T. Owens
2013. The hermeneutics of data and historical writing. *Writing history in the digital age*, 159.
- Gil, Y. and P. Groth
2011. LinkedDataLens: linked data as a network of networks. In *Proceedings of the sixth international conference on Knowledge capture*, Pp. 191–192. ACM.

- Granovetter, M.
1985. Economic action and social structure: The problem of embeddedness. *American journal of sociology*, 91(3):481–510.
- Groth, P. and Y. Gil
2011. Linked data for network science. In *Proceedings of the First International Conference on Linked Science-Volume 783*, Pp. 1–12. CEUR-WS. org.
- Hagberg, A. and D. Conway
2010. Hacking social networks using the python programming language. *Sunbelt 2010, Riva del Garda, Italy*.
- Hagberg, A., P. Swart, and D. S Chult
2008. Exploring network structure, dynamics, and function using networkx. Technical report, Los Alamos National Lab.(LANL), Los Alamos, NM (United States).
- Krech, D.
2006. Rdfliib: A python library for working with rdf.
- Lebo, T., S. Sahoo, D. McGuinness, K. Belhajjame, J. Cheney, D. Corsar, D. Garijo, S. Soiland-Reyes, S. Zednik, and J. Zhao
2013. Prov-o: The prov ontology. *W3C recommendation*.
- Meroño-Peñuela, A., A. Ashkpour, C. Guéret, and S. Schlobach
2015. CEDAR: the Dutch Historical Censuses as Linked Open Data. *Semantic Web*, 8(2):297–310.
- project Agents of Change: Women Editors, E. and .-. W. Socio-Cultural Transformation in Europe
2015. Wechanged.
- Schelstraete, J. and M. Van Remoortel
2019. Towards a sustainable and collaborative data model for periodical studies. *Media History*, 25(3):336–354.
- Segaran, T., C. Evans, and J. Taylor
2009. *Programming the Semantic Web: Build Flexible Applications with Graph Data*. "O'Reilly Media, Inc."
- Simmel, G.
1955. The web of group-affiliations. conflict and the web of groupaffiliations.
- Thornton, K., J. M. Birkholz, M. Van Remoortel, and S.-N. Kenneth
forthcoming. Working paper on bringing to light to women editors of the past through linked open data on wikidata.
- Thornton, K., E. Cochrane, T. Ledoux, B. Caron, and C. Wilson
2017. Modeling the domain of digital preservation in wikidata. In *Proceedings of ACM International Conference on Digital Preservation, Kyoto, Japan*.
- Van Remoortel, M.; Birkholz, J. S.-J. A. M. B. C. D. C. F. E. D. G.-G. N. G. M. J. A. V. G.
2020. Wechanged database.
- Van Steen, M.
2010. Graph theory and complex networks. *An introduction*, 144.

(W3C), W. W. W. C.

2011. W3c semantic web activity.

Wasserman, S., K. Faust, et al.

1994. *Social network analysis: Methods and applications*, volume 8. Cambridge university press.