

Data-Transformation on historical data using the RDF Data Cube Vocabulary

Sebastian Bayerl, Michael Granitzer
Department of Media Computer Science
University of Passau

SWIB15 – Semantic Web in Libraries
22.10.2015

Overview

- Motivation
- Vocabulary and Dataset
- Problem Setting and Approach
- Workflow
- Contributions

Motivation

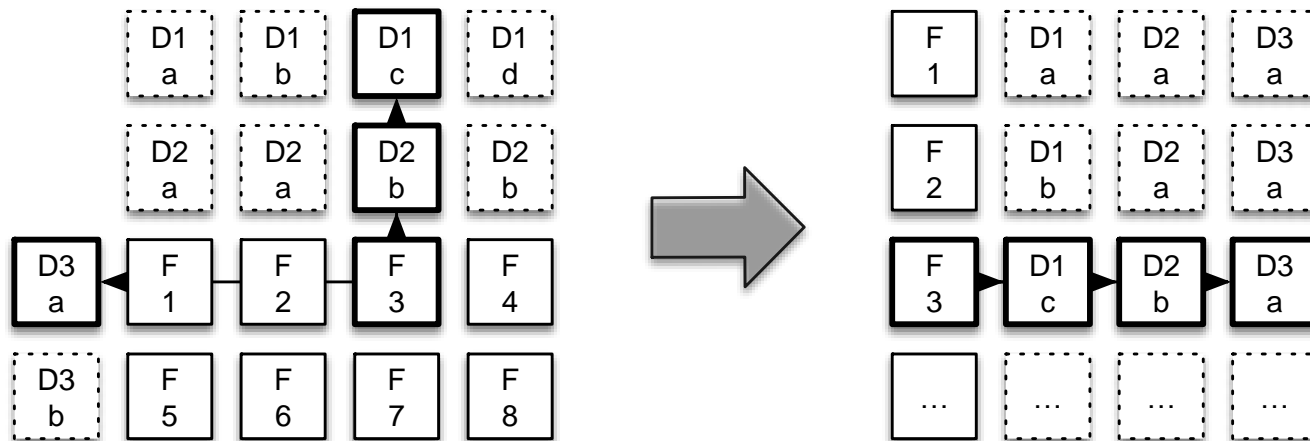
- Statistical and historical data source
 - Statistics of the German Reich (Digitalized)
- Access the encapsulated knowledge
- Data Analytics and Recommendation
 - Using Linked Data (RDF Data Cube Vocabulary)
- But first: Data Integration
 - Data Cleaning, -Transformation and -Fusion

A. Nachweisung über die Abfertigungen in Bezug auf die

[illegible]

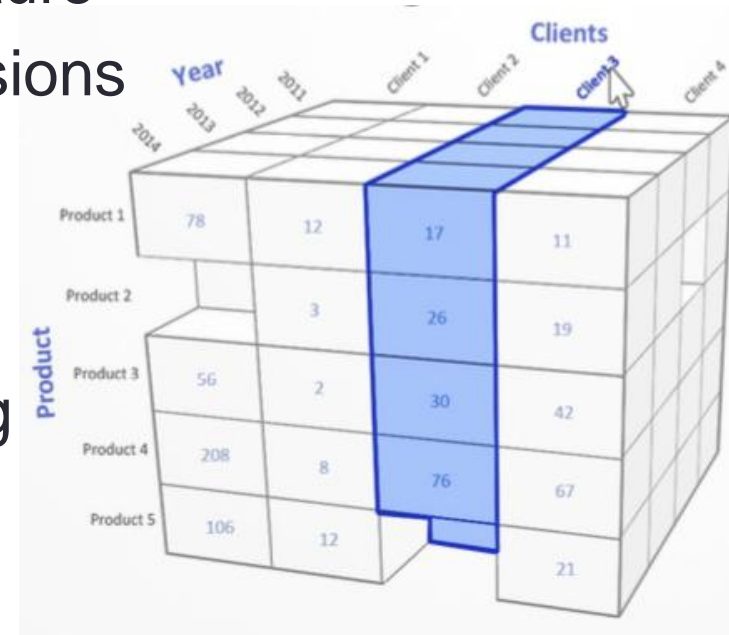
Target structure

-1	0	1	2	3	4	5
0	25330	Einfuhr	im 1. Quartal 1873.	Glasscherben, Bruchglas.	1. (132.) Pos. 1 a.	Preussen
1	21861	Einfuhr	im 1. Quartal 1872.	Glasscherben, Bruchglas.	1. (132.) Pos. 1 a.	Preussen
2	337	Einfuhr	im 1. Quartal 1873.	Glasscherben, Bruchglas.	1. (132.) Pos. 1 a.	Bayern
3	378	Einfuhr	im 1. Quartal 1872.	Glasscherben, Bruchglas.	1. (132.) Pos. 1 a.	Bayern
4	474	Einfuhr	im 1. Quartal 1873.	Glasscherben, Bruchglas.	1. (132.) Pos. 1 a.	Baden
5	120	Einfuhr	im 1. Quartal 1872.	Glasscherben, Bruchglas.	1. (132.) Pos. 1 a.	Baden



Data Cubes and OLAP

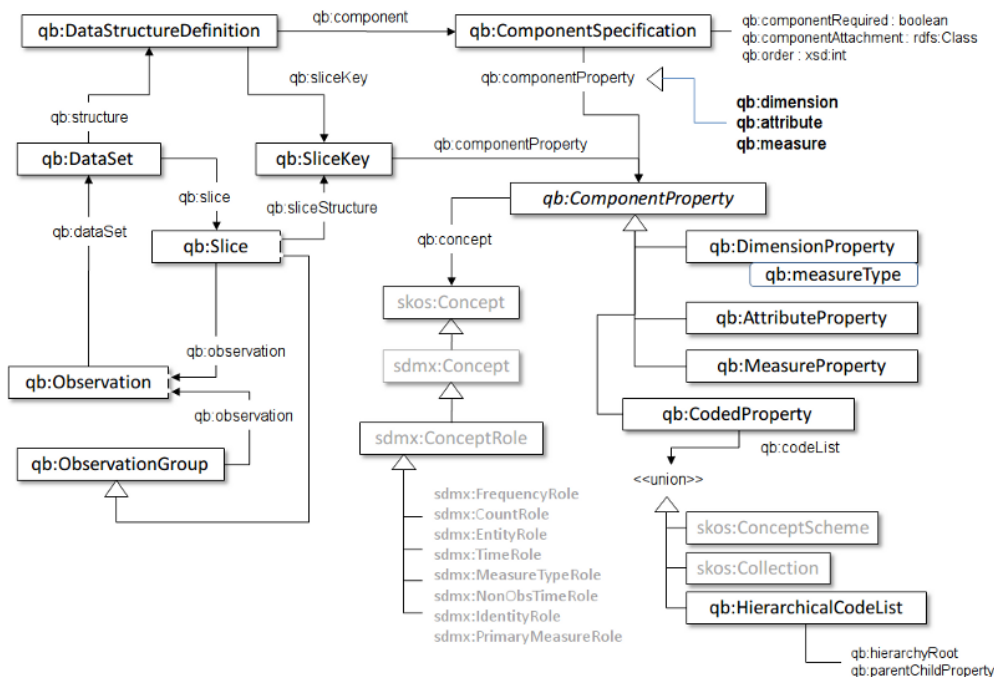
- Cube: Multi-dimensional data structure
- Observation: measures and dimensions
 - Measure: numerical fact
 - Dimension: describes the fact(s)
- Enables Data Analytics
- OLAP: Online Analytical Processing
 - Slicing, Dicing, Roll-Up,...



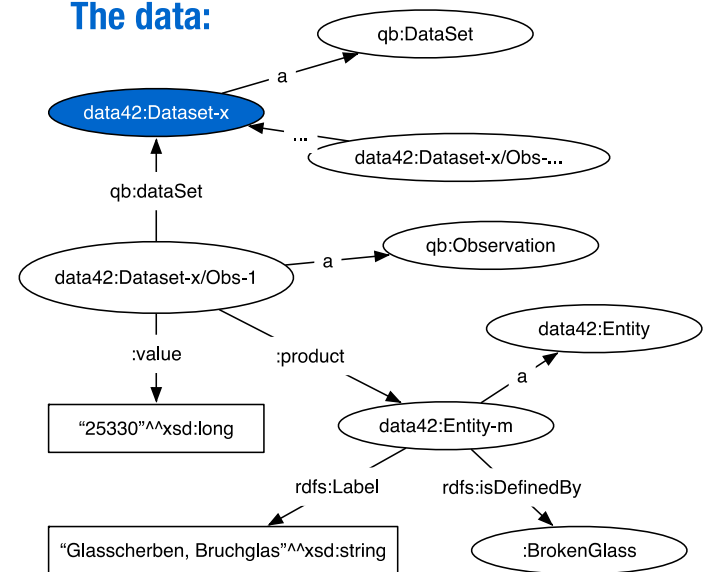
-1	0	1	2	3	4	5
0	25330	Einfuhr	im 1. Quartal 1873.	Glasscherben, Bruchglas.	1. (132.) Pos. 1 a.	Preussen
1	21861	Einfuhr	im 1. Quartal 1872.	Glasscherben, Bruchglas.	1. (132.) Pos. 1 a.	Preussen
2	337	Einfuhr	im 1. Quartal 1873.	Glasscherben, Bruchglas.	1. (132.) Pos. 1 a.	Bayern
3	378	Einfuhr	im 1. Quartal 1872.	Glasscherben, Bruchglas.	1. (132.) Pos. 1 a.	Bayern
4	474	Einfuhr	im 1. Quartal 1873.	Glasscherben, Bruchglas.	1. (132.) Pos. 1 a.	Baden
5	120	Einfuhr	im 1. Quartal 1872.	Glasscherben, Bruchglas.	1. (132.) Pos. 1 a.	Baden

The RDF Data Cube Vocabulary

- RDF based vocabulary
- Models an OLAP Data Cube
- Interlink components with existing concepts



The data:



Examples 1

III. Uebersicht der im Seeverkehr angekommenen und abgegangenen Schiffe nach den Flaggen und nach den Ländern (Küstenstrecken) der Herkunft und Bestimmung für das Jahr 1873.

Länder bezw. Küstenstrecken der Herkunft und Bestimmung.		A n g e k o m m e n					A b g e g a n g e n				
		Dampfschiffe mit schrägen Ziffern, in den Hauptzahlen mit enthalten.									
		Mit Ladung.		In Ballast oder leer.		Be- satzung.	Mit Ladung.		In Ballast oder leer.		Be- satzung.
		Schiffe.	Reg.-Tons.	Schiffe.	Reg.-Tons.		Schiffe.	Reg.-Tons.	Schiffe.	Reg.-Tons.	
1.		2.	3.	4.	5.	6.	7.	8.	9.	10.	11.
A. In den Preussischen Hafenplätzen.											
1. Deutsche Schiffe überhaupt.											
I. Deutsches Reich.											
Preussischer Staat.	Preussen	1 092	94 996	61	8 434	7 003	761	63 652	520	60 574	7 751
		291	40 324	13	1 970	3 601	220	29 959	35	5 856	3 446
	Pommern	2 277	153 380	617	55 706	14 973	1 858	132 474	358	17 185	11 886
		598	84 266	62	10 106	7 740	582	80 173	51	3 590	7 155
	Schleswig-Holstein an der Ostsee . . .	4 366	105 416	1 982	53 638	16 542	4 563	124 958	1 560	27 827	16 155
		632	43 825	82	6 373	4 504	667	49 302	68	3 003	4 862
	„ „ an der Nordsee . . .	984	22 763	337	8 054	2 859	921	22 070	481	11 219	3 030
		4	220	1	91	35	4	218	1	42	34
	Hannover, östl. Theil	286	7 934	116	2 663	1 054	368	13 271	141	3 374	1 427
		23	2 232	—	—	220	25	2 452	1	60	253
	westl. Theil einschl. des Jadegeb.	910	97 248	728	18 981	4 536	1 225	58 515	—	—	—

Seereisen Deutscher Schiffe im Jahre 1874.

Noch: III. A. Reisen zwischen ausserdeutschen Häfen, zusammengestellt nach den Küstenstrecken des Abgangs.

Länder bzw. Küstenstrecken des Abgangs und der Bestimmung.	Mit Ladung.		In Ballast oder leer.	
	Schiffe.	Reg.-Tons.	Schiffe.	Reg.-Tons.
1.	2.	3.	4.	5.

Noch: 7. Deutsche Schiffe überhaupt.

Noch: Von Grossbritannien und Irland nach:

Südamerika a. atl. Meere, südl. v. Brasilien	56	20033	—	—
Chile	12	6853	—	—
dem übrigen Südamerika am stillen Meere	13	7749	—	—
Egypten am mittelländischen Meere	11	5598	—	—
Kapland mit Natal	26	6888	—	—
Afrika am atlantischen Meere	29	8163	9	1038
" " indischen und rothen Meere	2	1116	—	—
Asien a. mittell. u. schwarz. Meere (Levante)	1	918	—	—
dem übrigen Vorderasien bis Ostindien	2	1538	—	—
Ostindien mit den indischen Inseln	35	25566	1	1030
China	11	5962	—	—
Australien mit den Inseln im stillen Meere	3	1278	1	1119
Ueberhaupt	1497	540494	711	208348

Von den Niederlanden nach:

dem Europ. Russland a. weiss. M. u. Eismeere	—	—	10	1838
" " " an der Ostsee	22	22925	37	6377
Schweden	8	1043	3	550

Problem Setting



- Data is encapsulated in multiple files
- Unusable for sophisticated Data Analysis
- Normalization of complex structured data
- Dirty and faulty data, structure or annotations
- Lots of similar problems in a huge dataset

Approach

- Use the RDF Data Cube Vocabulary
 - Enables: Interlinking, merging and analytics
- Use an incremental workflow
 - Identify fine-granular transformations
- Implement the research prototype with GUI
 - Select, configure and chain transformations (save/load)
 - HTML preview

The screenshot shows the Statistics2Cubes application window. The main area displays a table with 12 rows and 6 columns. The first column contains row numbers (3-12), the second column contains location names, and the remaining four columns contain numerical values. The table is titled 'label' and '2'. On the left, a sidebar lists transformations: 'Show original table', 'NormalizeCompoundTables', 'TrimValues', 'SanityNotEmpty', and 'CreateHeaders {Anzahl, Kateg}'. The 'NormalizeCompoundTables' transformation is selected. Below the sidebar is a 'Transform' button and a version indicator '1.1'. On the right, a log window shows the execution status of the transformations.

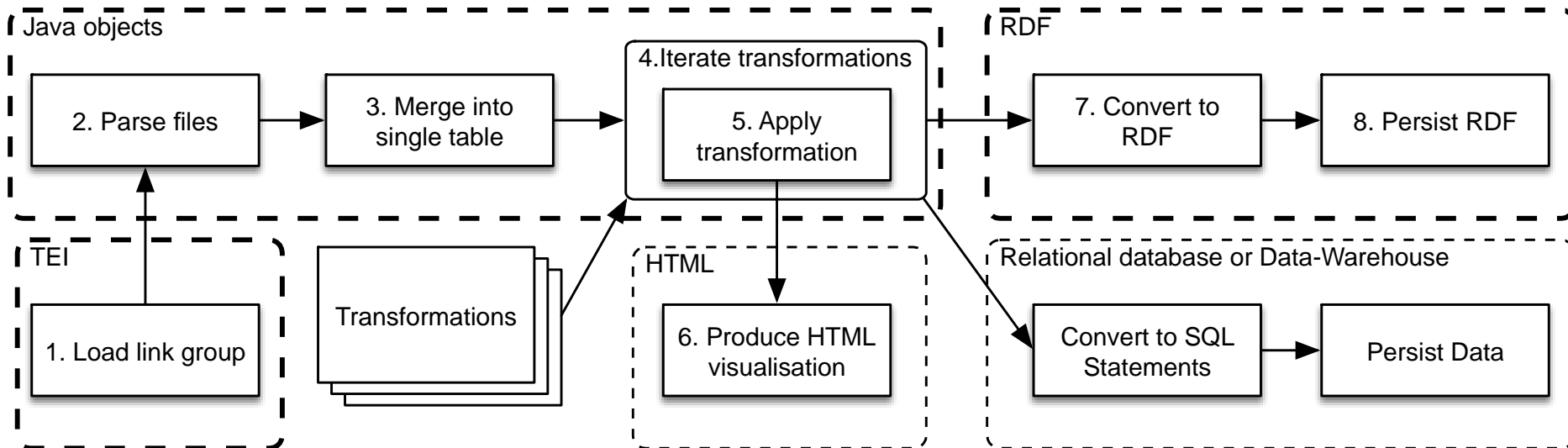
		Ctr.	Ctr.	Ctr.	Ctr.	Ctr.
3						
4	1. Baumwollengarn. (12 bis 14.)					
5	Königsberg i. Pr.	775	1442	2217	1464	753
6	Stettin	1328	1950	3278	2729	549
7	Elberfeld	742	6538	7280	3927	3353
8	Düsseldorf	1667	2343	4010	2557	1453
9	Leipzig	5442	6348	11790	6779	5011
10	Löbau	2521	6352	8873	5630	3243
11	Zittau	1144	3850	4994	3351	1643
12	Uebrige	1991	1279	3270	1779	1491

Log window output:

```

6 Table(s) loaded (and merged)
Table processed in 19 ms. Re
Table processed in 28 ms. Re
Table processed in 14 ms. Re
Transformations executed in
Transformation list valid.xml
  
```

Workflow



Transformations

1. Pre-Normalization

- Sanity checks
- Data Cleaning
- Fix structure (e.g. spans), data and annotations
- Delete row (e.g. repeating headers)
- 30 more...

2. Normalization

- Normalization
- Compound normalization: Horizontal or vertical partitions

3. Post-Normalization

- Add/merge/delete columns
- Add headers/disambiguation
- Add metadata
- ...

Advanced transformations

- Compound transformations
 - Combine multiple transformation
 - Fix more complex problems
 - E.g. find problematic cells and fix with existing transformation
- Transformation suggestions
 - Find common problems: Repeat symbol, annotation patterns
 - A step towards automation

Contributions

- Modular workflow for the Data Integration process
 - Definition of fine granular transformation steps
 - Reusable within the same or for other data sources
- Lift and enrich historical statistical data
 - Ready for Data Analytics
- Current dataset contains 32169 files
 - > 10% converted
 - 10 conversion chains



Thank you for your attention!



RDF Data Cube Vocabulary:

<http://www.w3.org/TR/vocab-data-cube/>



<https://github.com/bayerls/statistics2cubes>

Sebastian Bayerl

Department of Media Computer Science

University of Passau

bayerl@dimis.fim.uni-passau.de

BACKUP

Publication

- Bayerl, Sebastian, and Michael Granitzer. "Data-transformation on historical data using the RDF data cube vocabulary." Proceedings of the 15th International Conference on Knowledge Technologies and Data-driven Business. ACM, 2015.

Abstract

This work describes how XML-based TEI documents, containing statistical data, can be normalized, converted and enriched using the RDF Data Cube Vocabulary. In particular we focus on a statistical real world data set, namely the statistics of the German Reich around the year 1880, which are available in the TEI format. The data is embedded in complex structured tables, which are relatively easy to understand for humans but they are not suitable for automated processing and data analysis, without heavy pre-processing, due to their varying structural properties and differing table layouts. Therefore, the complex structured tables must be validated, modified and transformed, until they are suitable for the standardized multi-dimensional data structure - the data cube. This work especially focuses on the transformations necessary to normalize the structure of the tables. Performing validation- and cleaning-steps, resolving row- and column-spans and reordering slices are available transformations among multiple others. By combining exiting transformations, compound operators are implemented, which can handle specific and complex problems. The identification of structural similarities or properties can be used to automatically suggest sequences of transformations. A second focus is on the advantages, which come by using the RDF Data Cube Vocabulary. Also, a research prototype was implemented to execute the workflow and convert the statistical data into data cubes.