

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

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

Big Data Analytics: Big Data management and Data Integration
i-KNOW
22.10.2015

Overview

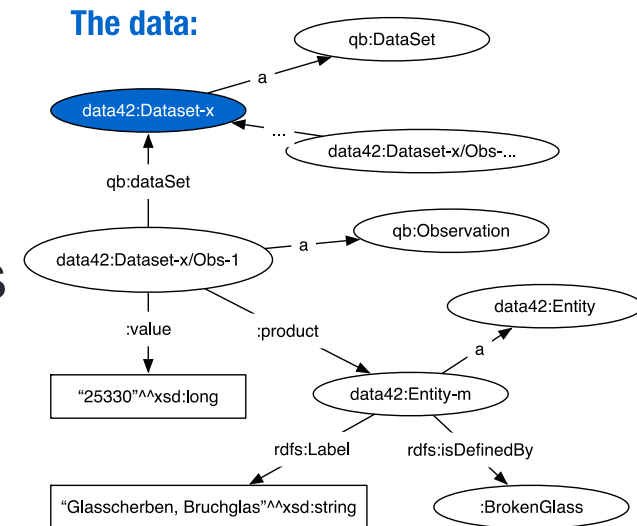
- Motivation
- Vocabulary and Dataset
- Problem Setting and Approach
- Workflow
- Contributions

Motivation

- EEXCESS Vision:
 - Unfold the treasure of cultural, educational and scientific long-tail content.
- Data Analytics!
 - Using Linked Data (RDF Data Cube Vocabulary)
 - Statistical and historical data source: German Reich statistics
- But first: Data Integration
 - Data Cleaning, -Transformation and -Fusion

The RDF Data Cube

- Cube: Multi-dimensional data structure
- Observation: measures and dimensions
 - Measure: numerical fact
 - Dimension: describes the fact(s)



-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

Example

Die deutschen überseeischen Auswanderer in den Jahren 1871 bis 1882.

Deutsche Auswanderer im Jahre	überhaupt (durch die amtliche deutsche Statistik nachweisbar)	davon wurden befördert													
		über				nach									
		Bremen.	Ham- burg.	Stet- tin.	Ant- werpen.	Amerika, nämlich						Afrika.	Asien.	Austra- lien.	
						den Ver- einigten Staaten v. A.	brit- tisch Nord- A.	Mexi- ko u. Zen- tral-A.	West- indien.	Brasi- lien.	ande- ren Thei- len v. A.				
1.	2.	3.	4.	5.	6.	7.	8.	9.	10.	11.	12.	13.	14.	15.	
1871	75 912	45 658	30 254	.	.	73 816	9	21	37	920	263	18	11	817	
1872	125 650	66 919	57 615	.	1 116	120 056	690	38	61	3 232	387	2	12	1 172	
1873	103 638	48 608	51 432	.	3 598	96 641	49	32	28	5 048	496	4	9	1 331	
1874	45 112	17 907	24 093	1 536	1 576	42 492	138	24	83	1 019	418	5	33	900	
1875	30 773	12 613	15 826	268	2 066	27 834	38	26	47	1 387	377	1	37	1 026	
1876	28 368	10 972	12 706	202	4 488	22 767	11	8	35	3 432	804	54	31	1 226	
1877	21 964	9 328	10 725	75	1 836	18 240	11	25	243	1 069	289	750	31	1 306	
1878	24 217	11 329	11 827	85	976	20 373	89	22	74	1 048	449	394	50	1 718	
1879	33 327	15 828	13 165	245	4 089	30 808	44	17	59	1 630	441	23	31	274	

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

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

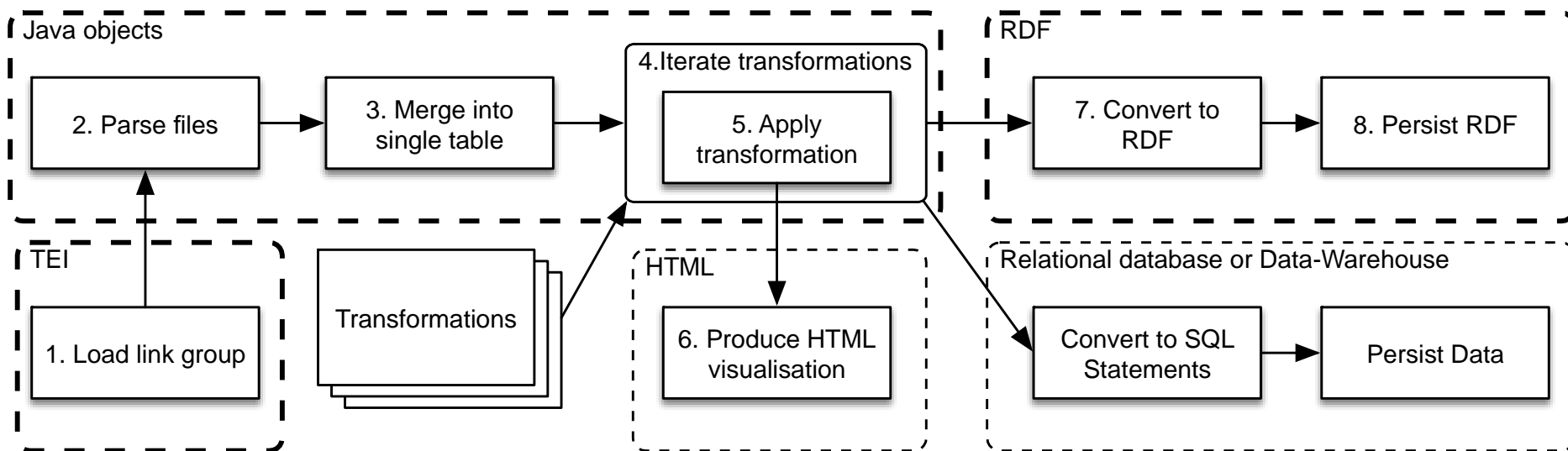
- Use the RDF Data Cube Vocabulary
 - Enables: Interlinking, merging and analytics
- Use an incremental workflow
 - Identify fine-granular transformations
- Implement the research prototype with GUI
 - Select, configure and chain transformations (save/load)
 - HTML preview

The screenshot shows the Statistics2Cubes application interface. At the top, there is a menu bar with 'File' and 'Help'. Below the menu bar, there is a dropdown menu set to 'NormalizeCompoundTables', a text input field containing 'label', and another text input field containing '2'. A 'Transform' button is visible below the input fields. The main area displays a table with 12 rows and 6 columns. The first row has headers 'Ctr.', 'Ctr.', 'Ctr.', 'Ctr.', and 'Ctr.'. The subsequent rows contain data for various locations. A sidebar on the left lists transformations: 'Show original table', 'NormalizeCompoundTables', 'TrimValues', 'SanityNotEmpty', and 'CreateHeaders {Anzahl, Kateg}'. At the bottom left, there is a version indicator '1.1' and an 'Export' button. On the right side, a log window displays the following text:

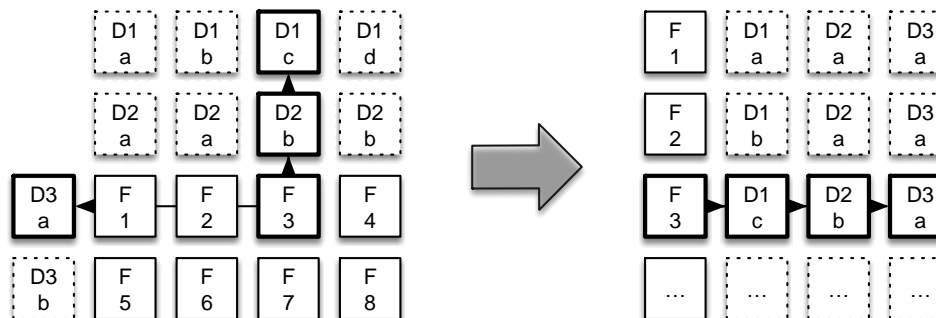
```
6 Table(s) loaded (and merged)
Table processed in 19 ms. Re
Table processed in 28 ms. Re
Table processed in 14 ms. Re
Transformations executed in
Transformation list valid.xml l
```

		Ctr.	Ctr.	Ctr.	Ctr.	Ctr.
3						
4	I. 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

Workflow



Normalize:



Transformations

Pre-Norm.

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

Normalization

- (Compound) normalization

Post-Norm.

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

- > 35 transformations
- Lots of Data Cleaning
- Generic- and structure dependent implementations

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 publication and Data Analytics
- Current dataset contains 32169 files
 - > 10% converted
 - 10 conversion chains



Thank you for your attention!

Question?



Contributions

- Define a modular and extendible conversion workflow
- Implemented as an open source research prototype
- Apply Semantic Web Standard to historical statistical data
- Enrich and lift the statistics of the German Reich

GUI

Statistics2Cubes

File Help

NormalizeCompoundTables label 2 ✓

Show original table

- NormalizeCompoundTables
- TrimValues
- SanityNotEmpty
- CreateHeaders {Anzahl, Kateg

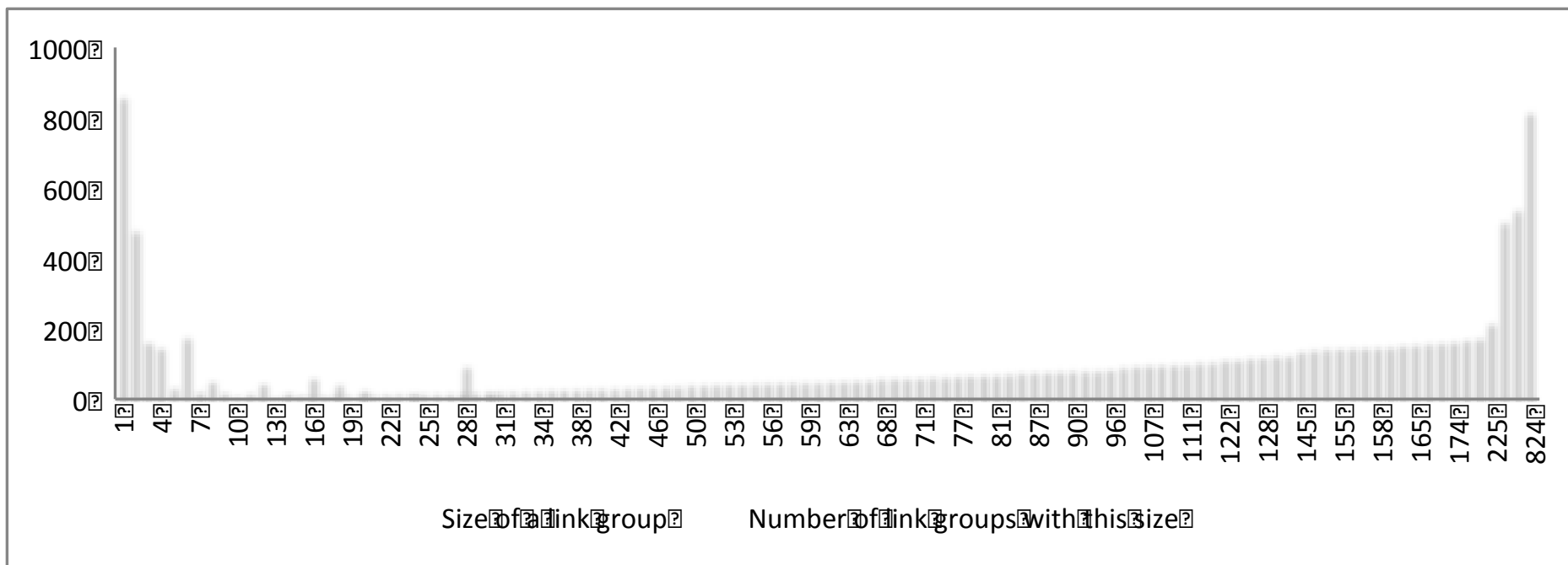
Transform

1.1 Export

3		Ctr.	Ctr.	Ctr.	Ctr.	Ctr.
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

6 Table(s) loaded (and merged)
 Table processed in 19 ms. Res
 Table processed in 28 ms. Res
 Table processed in 14 ms. Res
 Transformations executed in 8:
 Transformation list valid.xml lo

Backup



Example 3

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 Jadeb.	910	97 248	728	18 001	5 530	1 000	50 511	—	—	—	

The dataset

- Statistics of the German Reich around the year 1880
- Provided by the ZBW
- Data is available as TEI files
- One page = one table → Link groups
- Highly complex structured
- Hard to interpret, even for humans
- Unclean or even invalid structure or data

-1	0	1	2	3	4
0	I. In den freien Verkehr des Deutschen Zollgebiets getretene Waaren.				
1	Nummern der Waarenverzeichnisse resp. Tarifposition*)	Gebietstheile.	Einfuhr		
2			im 1. Quartal 1873.	im 1. Quartal 1872.	im 1. Quartal 1873 mehr (+) weniger (-).
3	1. (132.) Pos. 1 a.	Glasscherben, Bruchglas. — Ctr. brutto, zollfrei.			
4		Preussen	25330	21861	+3469
5		Bayern	337	378	-41
6		Baden	474	120	+354
7		Oldenburg	2252	2746	-494
8		Elsass-Lothringen	3106	1527	+1579
9		Uebrigtes Zollgebiet	204	256	-52