

An Evaluation of RDF Constraint Types to Validate Metadata on Person-Level and Aggregated Data

Thomas Bosch¹, Benjamin Zapilko¹, Joachim Wackerow¹, and Kai Eckert²

¹ GESIS – Leibniz Institute for the Social Sciences, Germany
{firstname.lastname}@gesis.org,

² University of Mannheim, Germany
kai@informatik.uni-mannheim.de

Abstract. ...

Keywords: RDF Validation, RDF Constraints, DDI-RDF Discovery Vocabulary, Disco, RDF Data Cube Vocabulary, Linked Data, Semantic Web

1 Introduction

2 Evaluation

2.1 Legend

Symbol	Description
X	Validation Successful
$> X$	Poor Performance/Scaling
\times	Very Poor Performance/Scaling
\dagger	SPARQL Endpoint Error

Table 1: Legend

Validation Successful: X indicates the number of raised constraint violation triples.

Poor Performance/Scaling: The performance of the implementation of the underlying SPARQL CONSTRUCT query is too poor to get all resulting constraint violation triples. Therefore, a limit of X result constraint violation triples is set. It is likely that there are more than X constraint violations. Although, the result set contains not the whole set of raised constraint violation triples, the constraint can be used as an indicator if there data not conforming to the constraint and to resolve constraint violations step by step. As part of future work, the performance will be improved.

Very Poor Performance/Scaling: The performance of the implementation of the underlying SPARQL CONSTRUCT query is too poor to get any results, even though a limit of result constraint violation triples is set. As part of future work, the performance will be improved.

3 Disco

– DwB Discovery Portal

- Metadata records for each study are available on the study landing pages.
- <http://dwb-dev.nsd.uib.no/portal/#/study/http://portal.dwbproject.org/resource/Study/de.gesis/zacat/ZA5597/en>
- Search interface: <http://dwb-dev.nsd.uib.no/portal/#/search?q=smoking&p=1&s=score&d=desc&f=text>

Abbr.	Disco Data Sets
Missy	Microdata Information System ³
DwB	DwB Discovery Portal ⁴
DDA-SND	DDA and SND DDI-RDF ⁵

Table 2: Disco Data Sets Abbreviations

³ <http://www.gesis.org/missy/eu/missy-home>

⁴ <http://dwb-dev.nsd.uib.no/portal>

⁵ <http://ddi-rdf.borsna.se/>

Data Sets	Counts									
	triples	disco:StudyGroup	disco:Study	disco:LogicalDataSet	disco:Universe	disco:Variable	disco:Question	disco:SummaryStatistics	disco:CategoryStatistics	skos:Concept
Missy										
DwB	2,332,802	0	1,387	1,367	2,796	446,806	0	0	0	0
DDA-SND	2,271,415	0	1,490	0	10,188	80,070	139,237	0	0	290,963

Table 3: Disco Data Sets Overview

Data Sets	SPARQL Endpoint
Missy	XXXXX
DwB	http://dwb-dev.nsd.uib.no/sparql
DDA-SND	http://ddi-rdf.borsna.se/endpoint/

Table 4: Disco SPARQL Endpoints

4 Data Cube

- overview over Data Cube data sets⁶
- <http://ontologycentral.com/>

⁶ <http://270a.info/>; <http://datahub.io/de/dataset?tags=format-qb>

Abbr.	Data Cube Data Sets
ECB	European Central Bank ⁷
UIS	UNESCO Institute for Statistics ⁸
IMF	International Monetary Fund ⁹
BFS	Bundesamt für Statistik - Swiss Federal Statistics ¹⁰
FAO	Food and Agriculture Organization of the United Nations ¹¹
WB	World Bank ¹²
FRB	Federal Reserve Board ¹³
TI	Transparency International ¹⁴
OECD	Organisation for Economic Co-operation and Development ¹⁵
BIS	Bank for International Settlements ¹⁶
ABS	Australian Bureau of Statistics ¹⁷
IEEE-VIS	IEEE VIS Source Data
ACORN-SAT	Australian Climate Observations Reference Network - Surface Air Temperature Dataset
HDP	HealthData.gov Platform (HDP) on the Semantic Web
Eurostat	The Eurostat Linked Data
Asturias	Nomenclator Asturias
ISTAT	ISTAT Immigration (LinkedOpenData.it)
ICANE	Statistical Office of Cantabria (Instituto Cántabro de Estadística, ICANE)
EE-2009	European Election Results 2009
EU-B	Standard Eurobarometer
ECB-S	European Central Bank Statistics (PublicData.eu)
CPV-2008	Common Procurement Vocabulary (CPV) 2008
CPV-2003	Common Procurement Vocabulary (CPV) 2003
GHI-2011	Global Hunger Index (GHI) 2011

Table 5: Data Cube Data Sets Abbreviations

⁷ <http://www.ecb.europa.eu/home/html/index.en.html>

⁸ <http://www.uis.unesco.org/Pages/default.aspx>

⁹ <http://www.imf.org/external/index.htm>

¹⁰ <http://www.bfs.admin.ch/>

¹¹ <http://www.fao.org/home/en/>

¹² <http://www.worldbank.org/>

¹³ <http://www.federalreserve.gov/>

¹⁴ <http://www.transparency.org/>

¹⁵ <http://www.oecd.org/>

¹⁶ <http://www.bis.org/>

¹⁷ <http://abs.gov.au/>

Data Sets	Counts				
	triples	qb:DataSet	qb:DataStructureDefinition	qb:Observation	qb:Slice
ECB	468,899,474	55	46	>11,000,000	428,698
UIS	10,400,534	5	5	1,437,651	0
IMF	35,688,446	4	8	3,603,719	0
BFS	1,533,743	0	0	8	0
FAO	53,000,000	10	10	>7,100,000	0
WB	174,006,552	9,466	59	>17,000,000	0
FRB	185,266,900	49	98	>9,500,000	0
TI	52,233	6	6	3,928	0
OECD	304,995,160	136	140	>12,000,000	0
BIS	54,197,482	6	12	3,606,466	47,914
ABS	2,357,400,000	253	257	>11,000,000	0
IEEE-VIS	19,935,340	0	0	1,350	0
ACORN-SAT	61,406,503				
HDP	1,934,046,908				
Eurostat	8,000,000,000				
Asturias	4,508,050				
ISTAT	3,024,396				
ICANE	52,412				
EE-2009	3,165				
EU-B	1,193,494				
ECB-S	10,000,000				
CPV-2008	803,311				
CPV-2003	546,135				
GHI-2011	19,369 (dump)				

Table 6: Data Cube Data Sets Overview

Data Sets	SPARQL Endpoints
ECB	http://ecb.270a.info/sparql
UIS	http://uis.270a.info/sparql
IMF	http://imf.270a.info/sparql
BFS	http://bfs.270a.info/sparql
FAO	http://fao.270a.info/sparql
WB	http://worldbank.270a.info/sparql
FRB	http://frb.270a.info/sparql
TI	http://transparency.270a.info/sparql
OECD	http://oecd.270a.info/sparql
BIS	http://bis.270a.info/sparql
ABS	http://abs.270a.info/sparql

Table 7: Data Cube Data Sets SPARQL Endpoints

Data Model Consistency	Data Sets											
	ECB	UIS	IMF	BFS	FAO	WB	FRB	TI	OECD	BIS	ABS	IEEE-VIS
DATA-MODEL-CONSISTENCY-01	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
DATA-MODEL-CONSISTENCY-02	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
DATA-MODEL-CONSISTENCY-03	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
DATA-MODEL-CONSISTENCY-04	✗	✓	✓	✓				✓				
DATA-MODEL-CONSISTENCY-05												✓
DATA-MODEL-CONSISTENCY-06												
DATA-MODEL-CONSISTENCY-07												
DATA-MODEL-CONSISTENCY-08												
DATA-MODEL-CONSISTENCY-09												
DATA-MODEL-CONSISTENCY-10												
DATA-MODEL-CONSISTENCY-11												

Table 8: Evaluation of Data Cube Data Sets

	Data Sets											
	ECB	UIS	IMF	BFS	FAO	WB	FRB	TI	OECD	BIS	ABS	IEEE-VIS
Existential Quantifications												
EXISTENTIAL-QUANTIFICATIONS-01	9	✓	11	7	8	77	8	9	7	8	7	✓
EXISTENTIAL-QUANTIFICATIONS-02	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
EXISTENTIAL-QUANTIFICATIONS-03	✓	✓	✓	✓	✓	59	✓	6	✓	✓	✓	✓
EXISTENTIAL-QUANTIFICATIONS-04	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓

Table 9: Evaluation of Data Cube Data Sets

	Data Sets											
	ECB	UIS	IMF	BFS	FAO	WB	FRB	TI	OECD	BIS	ABS	IEEE-VIS
Cardinality Restrictions												
MINIMUM-QUALIFIED-CARDINALITY-RESTRICTIONS-01												
EXACT-UNQUALIFIED-CARDINALITY-RESTRICTIONS-01	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
EXACT-QUALIFIED-CARDINALITY-RESTRICTIONS-01	✗	118	8	8	✗	✗	✗	✓	✗	12	✗	1,350
EXACT-QUALIFIED-CARDINALITY-RESTRICTIONS-02	✓	✓	✓	✓	✓	1	✓	✓	✓	✓	✓	✓

Table 10: Evaluation of Data Cube Data Sets

	Data Sets											
	ECB	UIS	IMF	BFS	FAO	WB	FRB	TI	OECD	BIS	ABS	IEEE-VIS
Structure												
STRUCTURE-01	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
STRUCTURE-02	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓

Table 11: Evaluation of Data Cube Data Sets

Constraints	Data Sets											
	ECB	UIS	IMF	BFS	FAO	WB	FRB	TI	OECD	BIS	ABS	IEEE-VIS
PROPERTY-DOMAIN-01												
PROPERTY-RANGES-01												
DISJOINT-PROPERTIES-01												
DISJOINT-CLASSES-01												
EQUIVALENT-PROPERTIES-01												
UNIVERSAL-QUANTIFICATIONS-01												
MEMBERSHIP-IN-CONTROLLED-VOCABULARIES-01												
CONTEXT-SPECIFIC-VALID-CLASSES-01												
CONTEXT-SPECIFIC-VALID-PROPERTIES-01												
RECOMMENDED-PROPERTIES-01												
VALUE-IS-VALID-FOR-DATATYPE-01												
VOCABULARY-01												

Table 12: Evaluation of Data Cube Data Sets

5 SKOS

- overview over SKOS data sets¹⁸
- overview over thesauri¹⁹ (table 14)

¹⁸ <http://datahub.io/de/dataset?tags=format-skos>

¹⁹ <http://datahub.io/de/dataset?tags=thesaurus>

Abbr.	Thesauri
TheSoz	Thesaurus for the Social Sciences ²⁰
STW	Thesaurus for Economics ²¹
AGROVOC	AGROVOC Multilingual agricultural thesaurus ²²
UNESCO	UNESCO Thesaurus ²³
TGN	The Getty Thesaurus of Geographic Names ²⁴
EARTH	Environmental Applications Reference Thesaurus ²⁵
ODT	Open Data Thesaurus ²⁶
SLD	Spanish Linguistic Datasets ²⁷
SSWT	Social Semantic Web Thesaurus ²⁸
GBA-GU	Thesaurus of the Geological Survey of Austria (GBA) - Geology Unit ²⁹
GBA-GTS	Thesaurus of the Geological Survey of Austria (GBA) - Geologic Time Scale ³⁰
GBA-L	Thesaurus of the Geological Survey of Austria (GBA) - Lithology ³¹
GBA-LU	Thesaurus of the Geological Survey of Austria (GBA) - Lithotectonic Unit ³²
GEMET	GEneral Multilingual Environmental Thesaurus ³³
EuroVoc	EuroVoc ³⁴
CECCT	Clean Energy and Climate Change Thesaurus ³⁵

Table 13: Thesauri Abbreviations

²⁰ <http://www.ecb.europa.eu/home/html/index.en.html>

²¹ <http://zbw.eu/stw/versions/latest/about>

²² <http://202.45.139.84:10035/catalogs/fao/repositories/agrovoc>

²³ <http://skos.um.es/sparql/>

²⁴ <http://vocab.getty.edu/sparql>

²⁵ <http://linkeddata.ge.imati.cnr.it/resource/EARTH/>

²⁶ <http://vocabulary.semantic-web.at/PoolParty/wiki/OpenData>

²⁷ <http://linguistic.linkeddata.es>

²⁸ <http://vocabulary.semantic-web.at/PoolParty/wiki/semweb>

²⁹ <http://resource.geolba.ac.at/>

³⁰ <http://resource.geolba.ac.at/>

³¹ <http://resource.geolba.ac.at/>

³² <http://resource.geolba.ac.at/>

³³ <http://www.eionet.europa.eu/gemet/>

³⁴ <http://open-data.europa.eu/de/data/dataset/eurovoc>

³⁵ <http://data.reegle.info/thesaurus/guide>

Thesauri	Counts						
	triples	skos:ConceptScheme	sko:Concept	skos:broader	skos:narrower	skos:hasTopConcept	skos:inScheme
TheSoz	439,153	1	8,426	13,705	13,706	0	48,529
STW	221,668	1	13,468	13,732	13732	7	13,180
AGROVOC	6,080,477	1	32,310	33,507	33,507	25	32,310
UNESCO	288,346	9	26,714	20,028	20,028	607	32,009
TGN	16,112,321	8	2,898,775	0	0	0	1,453,767
EARTH	9,287,364						
ODT	3,290						
SLD	7,629,211						
SSWT	64,698						
GBA-GU	25,718						
GBA-GTS	7,875						
GBA-L	9,317						
GBA-LU	9,504						
GEMET	372,889,229						
EuroVoc	64,477,774						
CECCT	191,336						

Table 14: Thesauri Overview

Thesauri	SPARQL Endpoints
TheSoz	http://lod.gesis.org/thesoz/sparql
STW	http://zbw.eu/beta/sparql/stw/query
AGROVOC	http://202.45.139.84:10035/catalogs/fao/repositories/agrovoc
UNESCO	http://skos.um.es/sparql/
TGN	http://vocab.getty.edu/
EARTH	http://linkeddata.ge.imati.cnr.it:8890/sparql
ODT	http://vocabulary.semantic-web.at/PoolParty/sparql/OpenData
SLD	http://linguistic.linkeddata.es/sparql
SSWT	http://vocabulary.semantic-web.at/PoolParty/sparql/semweb
GBA-GU	http://resource.geolba.ac.at/PoolParty/sparql/GeologicUnit
GBA-GTS	http://resource.geolba.ac.at/PoolParty/sparql/GeologicTimeScale
GBA-L	http://resource.geolba.ac.at/PoolParty/sparql/lithology
GBA-LU	http://resource.geolba.ac.at/PoolParty/sparql/tectonicunit
GEMET	http://semantic.eea.europa.eu/sparql
EuroVoc	http://open-data.europa.eu/de/linked-data
CECCT	http://poolparty.reegle.info/PoolParty/sparql/glossary

Table 15: Thesauri SPARQL Endpoints

	Data Sets
	TheSoz STW AGROVOC TGN UNESCO ODT SSWT GBA-GU GBA-GTS GBA-L GBA-LU CECCT
Data Model Consistency	
DATA-MODEL-CONSISTENCY-01	
DATA-MODEL-CONSISTENCY-02	
DATA-MODEL-CONSISTENCY-03	

Table 16: Thesauri Evaluation

Labeling and Documentation	Data Sets											
	TheSoz	STW	AGROVOC	TGN	UNESCO	ODT	SSWT	GBA-GU	GBA-GTS	GBA-L	GBA-LU	CECCT
LABELING-AND-DOCUMENTATION-01	8,426	11,508	19,829	✓	36	1,475	6	2	✓	107	486	
LABELING-AND-DOCUMENTATION-02	>1	✗	>100	✓	✓	✓	✓	✓	✓	✓	✓	✓
LABELING-AND-DOCUMENTATION-03	✓	✓	1	✓	✓	✓	✓	1	✓	✓	1	✓
LABELING-AND-DOCUMENTATION-04												
LABELING-AND-DOCUMENTATION-05	✓	✓	✗	✓	✓	✓	✓	1	✓	✓	✓	3
LABELING-AND-DOCUMENTATION-06	>1	✗	✗	✓	✓	✓	✓	✓	✓	✓	✓	✓

Table 17: Thesauri Evaluation

Structure	Data Sets											
	TheSoz	STW	AGROVOC	TGN	UNESCO	ODT	SSWT	GBA-GU	GBA-GTS	GBA-L	GBA-LU	CECCT
STRUCTURE-01	1	1,074	✓	✓	5	1	✓	✓	✓	✓	✓	✓
STRUCTURE-02												
STRUCTURE-03	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
STRUCTURE-04	2,906	>1,000	726	✓	12	124	84	256	>54	22	2,422	
STRUCTURE-05	✓	✓	✓	✓	90	5,150	✓	✓	✓	✓	9,864	
STRUCTURE-06	2	✓	✓	✓	✓	✓	✓	✓	✓	✓	113	
STRUCTURE-07	8	1,074	✓	✓	✓	✓	✓	✓	✓	✓	✓	
STRUCTURE-08												
STRUCTURE-09	6,741	>3,297	135	✓	2	16	25	✓	✓	✓	81	
STRUCTURE-10	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	

Table 18: Thesauri Evaluation

			Data Sets											
Language	Tag	Cardinality	TheSoz	STW	AGROVOC	TGN	UNESCO	ODT	SSWT	GBA-GU	GBA-GTS	GBA-L	GBA-LU	CECCT
LANGUAGE-TAG-CARDINALITY-01	9,435	13,468	98,894	✓		541	10,147	5,117	2,061	1,742	2,272	15,550		
LANGUAGE-TAG-CARDINALITY-02	8,222	36,936	✗	✓		265	3,627	2,212	635	631	1,253	9,607		
LANGUAGE-TAG-CARDINALITY-03	8,222	✓	135	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
LANGUAGE-TAG-CARDINALITY-04	✓	476	✗	50	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓

Table 19: Thesauri Evaluation

Constraints	Data Sets											
	TheSoz	STW	AGROVOC	TGN	UNESCO	ODT	SSWT	GBA-GU	GBA-GTS	GBA-L	GBA-LU	CECCT
PROPERTY-DOMAIN-01												
PROPERTY-RANGES-01												
DISJOINT-PROPERTIES-01												
DISJOINT-PROPERTIES-02												
DISJOINT-CLASSES-01												
EQUIVALENT-PROPERTIES-01												
UNIVERSAL-QUANTIFICATIONS-01												
CONTEXT-SPECIFIC-VALID-CLASSES-01												
CONTEXT-SPECIFIC-VALID-PROPERTIES-01												
RECOMMENDED-PROPERTIES-01												
VOCABULARY-01												

Table 20: Thesauri Evaluation

6 XKOS

- Nomenclature d’activités française (NAF): French classification expressed in XKOS. the French refinement of the NACE, because it has explanatory notes.
- Nomenclature des Professions et Catégories Socioprofessionnelles (PCS): French classification expressed in XKOS.
- Nomenclature des catégories juridiques (CJ): French classification expressed in XKOS.
- ISIC: has explanatory notes too.
- ISCO

	Counts
Data Sets	triples
Nomenclature d’activités française (NAF) ³⁶	
Nomenclature des Professions et Catégories Socioprofessionnelles (PCS) ³⁷	
Nomenclature des catégories juridiques (CJ) ³⁸	
ISIC	
ISCO	

Table 21: Statistical Classifications Overview

References

³⁶ <http://rdf.insee.fr/codes/index.html>

³⁷ <http://rdf.insee.fr/codes/index.html>

³⁸ <http://rdf.insee.fr/codes/index.html>