



Sofia Ribeiro Manso de Abreu e Silva

Catalogue of Linked Data Cube Descriptions

Dissertação de Mestrado

Dissertation presented to the Programa de Pós-Graduação em Informática of the Departamento de Informática, PUC-Rio as partial fulfillment of the requirements for the degree of Mestre em Informática.

Advisor: Prof. Marco Antonio Casanova

Rio de Janeiro
June 2013



Sofia Ribeiro Manso de Abreu e Silva

Catalogue of Linked Data Cube Descriptions

Dissertation presented to the Programa de Pós-Graduação em Informática do Departamento de Informática do Centro Técnico e Científico da PUC-Rio, as partial fulfillment of the requirements for the degree of Mestre.

Prof. Marco Antonio Casanova
Advisor
Departamento de Informática — PUC-Rio

Prof. Antonio Luz Furtado
Departamento de Informática — PUC-Rio

Prof. Luiz André Portes Paes Leme
UFF

Prof. José Eugenio Leal
Coordinator of the Centro
Técnico Científico — PUC-Rio

Rio de Janeiro, June 28th, 2013

All rights reserved

Sofia Ribeiro Manso de Abreu e Silva

Graduated in Applied Mathematics and Computation from the University of Aveiro (Portugal) in 2007. She joined the Master in Informatics at Pontifical Catholic University of Rio de Janeiro (PUC-Rio) in 2011. Her current research area is related to Statistical Linked Data and Semantic Web.

Bibliographic data

Ribeiro Manso de Abreu e Silva, Sofia

Catalogue of Linked Data Cube Descriptions / Sofia Ribeiro Manso de Abreu e Silva; advisor: Marco Antonio Casanova. — 2013.

85 f. : il. (color); 30 cm

1. Dissertação (mestrado) - Pontifícia Universidade Católica do Rio de Janeiro, Departamento de Informática, 2013.

Inclui bibliografia.

1. Informática – Teses. 2. Dados Estatísticos. 3. Linked Data. 4. Arquitetura de Mediação. 5. Triplificação. 6. RDF. 7. Data Cube Vocabulary. 8. R2RML. I. Casanova, Marco Antonio. II. Pontifícia Universidade Católica do Rio de Janeiro. Departamento de Informática. III. Título.

CDD: 004

Acknowledgements

First, I would like to express my gratitude to my advisor Marco Antonio Casanova for the useful comments, remarks and engagement through my learning process in this master. I am very thankful for his sharing of knowledge, support, motivation and trust during all my work. His encouragement and guidance was very important during this dissertation.

I would like to acknowledge the professors and colleagues, for the teachings, encouragement, help and friendship.

I also would like to thank the employees of the Informatics Department of PUC-Rio for their help, kindness and support.

To CAPES for the financial aid granted.

To my family, my sincere thanks. In particular, I would like to thank to my aunt Isabel and to my uncle António. Without them this master would not have been possible.

A big thanks to my friends for always being there for me.

I would like to make a special reference to my boyfriend Steve for his love and patience during these difficult times.

Last but not least, my deepest gratitude to my parents, Alice and João, for their unconditional love. I want to thank them for always believing in me, for giving me the strength to carry on and to make this happen. I will be grateful forever for their love.

Abstract

Ribeiro Manso de Abreu e Silva, Sofia ; Casanova, Marco Antonio (Advisor). **Catalogue of Linked Data Cube Descriptions.** Rio de Janeiro, 2013. 85p. MSc. Dissertation — Departamento de Informática, Pontifícia Universidade Católica do Rio de Janeiro.

Statistical Data are considered one of the major sources of information and are essential in many fields as they can work as social and economic indicators. A statistical data set comprises a collection of observations made at some points of a logical space and is often organized as what is called a data cube. The proper definition of the data cubes, especially of their dimensions, helps processing the observations and, more importantly, helps combining observations from different data cubes. In this context, the Linked Data principles can be profitably applied to the definition of data cubes, in the sense that the principles offer a strategy to provide the missing semantics of the dimensions, including their values.

This dissertation first describes a mediation architecture to help describing and consuming statistical data, exposed as RDF triples, but stored in relational databases. One of the features of this architecture is the Catalogue of Linked Data Cube Descriptions, which is described in detail in the dissertation. This catalogue has a standardized description in RDF of each data cube actually stored in statistical (relational) databases. Therefore, the main discussion in this dissertation is how to represent the data cubes in RDF, i.e., how to map the database concepts to RDF in a way that makes it easy to query, analyze and reuse statistical data in the RDF format.

Keywords

Statistical Data; Linked Data; Mediation Architecture; Triplification; RDF; Data Cube Vocabulary; R2RML.

Resumo

Ribeiro Manso de Abreu e Silva, Sofia ; Casanova, Marco Antonio.

Catálogo de Descrições de Cubos de Dados Interligados.

Rio de Janeiro, 2013. 85p. Dissertação de Mestrado — Departamento de Informática, Pontifícia Universidade Católica do Rio de Janeiro.

Dados estatísticos são considerados uma das principais fontes de informação e são essenciais em muitos campos, uma vez que podem funcionar como indicadores sociais e económicos. Um conjunto de dados estatísticos compreende um conjunto de observações feitas em determinados pontos de um espaço lógico e é muitas vezes organizado como o que se chama de cubo de dados. A definição correta dos cubos de dados, especialmente das suas dimensões, ajuda a processar as observações e, mais importante, ajuda a combinar as observações de diferentes cubos de dados. Neste contexto, os princípios de Linked Data podem ser proveitosamente aplicados à definição de cubos de dados, no sentido de que os princípios oferecem uma estratégia para proporcionar a semântica ausentes das suas dimensões, incluindo os seus valores.

Esta dissertação descreve inicialmente uma arquitetura de mediação para ajudar a descrever e consumir dados estatísticos, expostos como triplas RDF, mas armazenados em bancos de dados relacionais. Uma das características desta mediação é o Catálogo de Descrições de Cubos de Dados Interligados, que vai ser descrito em detalhes na dissertação. Este catálogo contém uma descrição padronizada em RDF para cada cubo de dados, que está realmente armazenado em cada banco de dados (relacional). Portanto, a principal discussão nesta dissertação é sobre a forma de representar em RDF cubos representando dados estatísticos e armazenados em bancos de dados relacionais, ou seja, como mapear os conceitos de banco de dados para RDF de uma forma em que seja fácil consultar, analisar e reutilizar dados estatísticos no formato RDF.

Palavras-chave

Dados Estatísticos; Linked Data; Arquitetura de Mediação; Triplificação; RDF; Data Cube Vocabulary; R2RML.

Table of Contents

1 Introduction	11
1.1. Motivation	11
1.2. Overview of the Mediation Architecture	13
1.3. The Three Stages of Data Consumption	15
1.4. Contributions	16
1.5. Related Work	16
1.6. Organization	18
2 Data Cube Description in RDF	20
2.1. Linked Data	20
2.2. RDF	21
2.3. Formal Definition of Data Cubes	23
2.4. Data Cube Vocabulary	24
2.5. R2RML	31
2.6. Triplification of Data Cubes	33
2.7. Data Cube Explanation and Example	35
2.8. Summary of Chapter 2	37
3 Proposed Architecture for the Catalogue of Linked Data Cube Descriptions	38
3.1. Overview of the Proposed Architecture	38
3.2. Catalogue of Linked Data Cube Descriptions	39
3.3. SameAs Generator	43
3.4. Data Cubes Charger	44
3.5. Interfaces	53
3.6. Summary of Chapter 3	54
4 Implementation of the Catalogue of Linked Data Cube Descriptions	55
4.1. Overview of the Implementation	55
4.2. Interface	56
4.2.1. Overview of the Interfaces	56

4.2.2. External Interface	57
4.2.3. Internal Interface	65
4.3. Summary of Chapter 4	79
5 Conclusions and Future Work	80
6 References	81

List of Figures

Figure 1 - Overview of the Mediation Architecture (Ruback et al. 2013)	13
Figure 2 - Representation of an RDF Triple (Klyne et al. 2004)	22
Figure 3 - Outline of the Data Cube Vocabulary (Cyganiak et al. 2013)	25
Figure 4 - An overview of R2RML	32
Figure 5 - Relational Schema for Cube Devotees	36
Figure 6 - Overview of the Proposed Architecture	38
Figure 7 - Database Schema	45
Figure 8 - Architecture of Virtuoso Jena Provider (Virtuoso 2012a)	55
Figure 9 - RDF triples stored in the catalogue	56
Figure 10 - Cube retrieved for the keyword “foreigners”	57
Figure 11 - Message informing that there are no cubes matching the keyword	57
Figure 12 - Output of External Interface for the cube selected	58
Figure 13 - Output of Internal Interface for the cube selected	66

List of Tables

Table 1 – Prefixes and Namespaces	26
Table 2 - Definition of the Cube Devotees	37
Table 3 - Definition of the Cube Residents	40
Table 4 - Definition of the connection to the RDB <i>project1</i>	41
Table 5 - Mapping of Cube Residents	42
Table 6 - Mapping of Dimension Race	42
Table 7 - RDF Triple representing the owl:sameAs of the Class Country	44
Table 8 - RDF Triple of a resource from Class Country	44
Table 9 - Description in RDF related with the Cube Residents	47
Table 10 - Correspondence Assertions	48
Table 11 - R2RML Mappings generated from CCA2 and DCA1	48
Table 12 - R2RML Mappings generated from CCA3 and DCA2	49
Table 13 - R2RML Mappings generated from CCA4 and DCA3	49
Table 14 - R2RML Mappings generated from CCA5 and DCA4	50
Table 15 - R2RML Mappings generated from CCA6 and DCA5	50
Table 16 - R2RML Mappings generated from CCA1, OCA1, OCA2, OCA3, OCA4, OCA5 and DCA6	52
Table 17 - RDF Triples retrieved by the External Interface about the Cube Foreigners	65
Table 18 - RDF Triples retrieved by the Internal Interface about the Cube Foreigners	79