

Copyright is owned by the Author of the thesis. Permission is given for a copy to be downloaded by an individual for the purpose of research and private study only. The thesis may not be reproduced elsewhere without the permission of the Author.

Dependencies in Complex-value Databases

Sebastian Link

A dissertation presented in partial fulfilment of the requirements for the
degree of Doctor of Philosophy in Information Systems at Massey
University

Supervisor:
Prof. Dr. Klaus-Dieter Schewe

Co-Supervisor:
Prof. Dr. Mike Hendy

Internal Examiner:
Associate Prof. Dr. Sven Hartmann

New Zealand Examiner:
Prof. Dr. Robert Goldblatt

Overseas Examiner:
Prof. Dr. Joachim Biskup

Date of Examination:
15.12.2004

Abstract

The relational data model has been the dominant model in database design for more than three decades. It considers data to be stored in matrices where rows correspond to individuals, columns correspond to attributes, and every cell contains a single atomic value. However, today's database technology trends, e.g. spatial, genetic or web-based data, require extended data models. Within the last decade new, complex-value data models such as the higher-order entity-relationship model, object-oriented data models, semi-structured data models, and XML have evolved which allow cells to contain lists, sets, multisets, trees, matrices or even more complex type constructors, references to other cells (which lead to infinite structures), and null values (indicating missing, unknown or vague data).

Matrices as such allow the storage of inconsistent data, invalid in the semantic sense. As this is not acceptable, additional requirements called dependencies have to be formulated when designing a database. The correct specification and use of dependencies needs a sound mathematical basis. For the relational data model more than 90 different classes of dependencies have been defined and studied intensively. The major problems in dependency theory are the axiomatisability of classes of dependencies, determination of the closure of a chosen set of dependencies (as certain dependencies can be implied by others) and the characterisation of semantically desirable properties for well-designed databases (such as absence of redundancies or abnormal update behavior) by syntactic properties on closed sets of dependencies.

With few exceptions research has only dealt with dependencies for the relational data model. Only recently, the emergence of XML as the standard format for web-based data and the rapidly increasing usage of persistent XML databases revealed the lack of a sound mathematical basis for complex-value data models. If they are expected to serve as first class data models they require a theoretical investigation of issues like integrity, consistency, data independence, recovery, redundancy, access rights, views and integration. The goal of this thesis is to develop a dependency theory for complex-value databases that is independent from any individual data model. Therefore, an abstract algebraic approach is taken that can be adapted to the presence of different combinations of type constructors such as records, lists, sets and multisets. Data models are classified according to the data types they support. In this framework the major objective is to initiate research on the following problems

- investigate the axiomatisation of important dependency classes, relevant to complex-value data models, by sound and complete sets of inference rules that permit the determination of all dependencies implied by some chosen set of dependencies.
- characterise semantically desirable properties by normal forms for complex-value data models and investigate whether these normal forms can always be achieved without violating other desirable properties.
- develop efficient algorithms for determining the closure of a chosen set of dependencies and for restructuring databases such that normal forms are satisfied and no information is lost.

In a single thesis it is impossible to consider all classes of relational dependencies in all different combinations of type constructors. Therefore the focus is put on extending two popular classes of relational dependencies: functional and multi-valued dependencies. The axiomatisation and implication of functional dependencies is investigated for all combinations of record, list, set and multiset type. Furthermore, a normal form with respect to functional dependencies in the presence of records and lists is proposed and semantically justified. It is also shown how to obtain databases which are in this normal form. Finally, axiomatisation and implication for the class of multi-valued dependencies and the class of functional and multi-valued dependencies are studied in the context of records and lists. The work of this thesis may lead to a unified dependency theory for complex-value data models.

Acknowledgement

I would like to thank *Klaus-Dieter Schewe* for giving me the opportunity to pursue an academic career under his supervision. Klaus-Dieter introduced me to theoretical computer science, and in particular to the topic of databases. His enthusiastic attitude towards research has motivated me since the beginning of my studies. Klaus-Dieter suggested to look at dependencies in the Higher-Order Entity-Relationship model, and it was later on that the investigations resulted in a more general treatment. The idea of classifying data models according to the data types they support is similar to the idea from [243] where object-oriented data models are classified according to different type systems. I also like to thank Klaus-Dieter for showing me how to work scientifically and for supporting decisions to reduce my teaching workload at critical times.

My thank also goes to *Mike Hendy* who kindly agreed to act as co-supervisor of this thesis.

Special thanks go to *Sven Hartmann* for having numerous fruitful discussions on the subject of this work. Sven made a lot of suggestions for improving the outcome, in particular on the mathematical rigorousness of various notions and proof arguments. Sven has always been prepared to talk about any issues, even in times he was very busy. I have learned a great deal from him.

I would further like to acknowledge *Bernhard Thalheim* and *Joachim Biskup* whose great experiences in the field of databases I have benefited from a lot.

Finally, I would like to thank my parents *Karla* and *Hans-Jürgen Link* who have always been there for me and supported me in every way possible.

Moreover, I am very thankful to my partner *Toni Floyd* who can always solve all my problems at once just by being who she is.

I dedicate this thesis to my parents,
Karla and Hans-Jürgen Link.

Table of Contents

1 Introduction	5
1.1 Relational Dependency Theory	5
1.1.1 Relational Dependencies	7
1.1.2 Functional Dependencies	8
Axiomatisation.	9
Implication Problem.	10
Boyce-Codd Normal Form.	11
1.1.3 Multi-valued Dependencies	13
Axiomatisation.	13
Implication Problem.	14
Minimality and Complementation Rule.	15
1.1.4 Additional Remarks and Literature	15
1.2 Challenges with Complex-value Databases	17
1.2.1 Extensions to the Relational Data Model	17
Semantic Data Models.	18
The Nested Relational Data Model.	19
Object-Oriented and Object-Relational Data Models.	20
Hypertext Datamodels.	21
Complex Objects in other Fields of Application.	22
1.2.2 Real-World Examples for Complex Constraints	24
Bioinformatics.	24
Image Processing.	28
Retailers.	29
1.3 Contributions	30
1.4 Outline	32
2 The Algebra of Nested Attributes	35
2.1 Brouwerian Algebras	35
2.2 Nested Attributes	38
2.2.1 Subattributes	40
2.2.2 The Brouwerian algebra of Subattributes	41
2.2.3 Notation, Examples and Intuition	44
2.3 Formalisation of Real-World Examples	48

2.4 Brouwerian algebras in the Literature	49
3 Functional Dependencies in the Presence of Lists	51
3.1 Axiomatisation	52
3.1.1 Definition of FDs	52
3.1.2 Implication and Derivation	53
3.1.3 The generalised Armstrong Axioms	56
3.1.4 Completeness	59
3.1.5 Dependencies for Keys	62
3.2 Implication Problem	63
3.2.1 The Closure	63
3.2.2 A first Approach	64
3.2.3 A different Perspective	66
3.2.4 A linear time Algorithm	68
3.2.5 Applications	71
3.3 Nested List Normal Form	72
3.3.1 Trivial FDs	73
3.3.2 The Notion of Redundancy	73
3.3.3 Boyce-Codd and Nested List Normal Form	76
3.3.4 NLNF - The same fact is only stored once	77
3.3.5 Characterising NLNF	78
3.3.6 Update Anomalies	80
3.4 Decomposition into NLNF	85
3.4.1 FDs and Decompositions	85
3.4.2 The Decomposition Algorithm	86
3.4.3 Problems with NLNF decomposition	90
4 Functional and Multi-valued Dependencies in the Presence of Lists	93
4.1 Axiomatisation	94
4.1.1 Definition and First Results	94
4.1.2 Trivial MVDs	96
4.1.3 MVDs are Binary Join Dependencies	96
4.1.4 Sound Inference Rules	97
4.1.5 Dependency Basis	103
4.1.6 Completeness	104
4.2 Minimality	108
4.3 Brouwerian-Complement Rule	114
4.4 Implication Problem	117
4.4.1 The Algorithm	117
4.4.2 Correctness	120
4.4.3 Complexity	128
4.4.4 Applications	130
4.5 The Class of Multi-valued Dependencies	131

4.5.1	Axiomatisation	131
4.5.2	Minimality	133
4.5.3	Implication Problem	135
4.5.4	A different Perspective for MVDs	136
4.6	Related and Future Work	137
5	Functional Dependencies in the Presence of Lists, Sets and Multisets	139
5.1	Axiomatisation	140
5.1.1	The Failure of the Extension Rule	140
5.1.2	Reconcilable Attributes	141
5.1.3	Soundness and some useful Inference Rules	143
5.1.4	Completeness	144
	Technical Lemmata.	145
	The Case of Sets.	147
	The Case of Multisets.	148
	The Main Lemma.	152
	The Main Theorem.	153
5.1.5	A Note on Reconcilability	154
5.2	Minimality	155
5.3	Minimal Axiomatisations for all Combinations	158
5.4	Implication Problem	159
5.4.1	The Closure	159
5.4.2	Units of Nested Attributes	160
5.4.3	Computing the Closure	163
5.4.4	Correctness	164
5.4.5	Complexity	167
5.4.6	Some Applications	169
5.5	The Implication Problem for all Combinations	170
5.6	Related Work	170
6	Summary	176
6.1	Main Results	176
6.2	Open Problems	178

List of Figures

1.1	Research on Dependency Theory	31
1.2	The Boolean algebra of Type Constructors	34
2.1	A Brouwerian algebra that is not an algebra of any nested attribute.	44
2.2	The Brouwerian algebra \mathcal{B} of $K\{L(A, M[N(B, C)])\}$	45
2.3	The poset (J, \leq) of the join-irreducible elements of \mathcal{B}	45
2.4	Brouwerian algebra of closed subsets of PO -space on (J, \leq)	46
2.5	The Brouwerian algebra of $K\{M(O\{A\}, P\{B\})\}$	47
2.6	Mathematical Concepts and Physics.	49
3.1	NLNF decomposition Tree of Example 3.22.	88
3.2	NLNF decomposition Tree of Example 3.23.	90
4.1	The Boolean algebra of $L(A)^M$	104
4.2	The subattribute basis of $K[L(M[N(A, B)], C)]$	105
4.3	Initialisation for $DepB_{\text{alg}}(X)$	119
4.4	$DepB_{\text{alg}}(X)$ after its first Update.	120
4.5	$DepB_{\text{alg}}(X)$ after its second Update.	120
4.6	Final State for $DepB_{\text{alg}}(X)$ from Example 4.9.	121
5.1	Identifying Terms of the Algebra $K\{L(A, M[N(B, C)])\}$	145
5.2	The closure \mathcal{X}^+ of $\mathcal{X} = K\{L(A)\}$	148
5.3	Illustration of Lemma 5.12	149
5.4	Illustration of Lemma 5.13	151
5.5	The structure of $M = K(J[A], O\{P(B, Q\{C\})\})$	152
5.6	The structure of Subalgebras in Example 5.4	152
5.7	Upper Complexity Bounds for the Implication Problem in the Presence of various Types	170
5.8	An XML data tree carrying some functional dependency.	173
5.9	An XML document corresponding to the XML data tree in Fig. 5.8.	173
5.10	Another XML data tree still carrying some functional dependency.	174
6.1	Subattribute Lattice of a Union-valued Attribute	181