

Encyclopedia of Data Warehousing and Mining

Second Edition

John Wang
Montclair State University, USA

Volume II
Data Pro-I

Information Science
REFERENCE

INFORMATION SCIENCE REFERENCE

Hershey • New York

Director of Editorial Content: Kristin Klinger
Director of Production: Jennifer Neidig
Managing Editor: Jamie Snavelly
Assistant Managing Editor: Carole Coulson
Typesetter: Amanda Appicello, Jeff Ash, Mike Brehem, Carole Coulson, Elizabeth Duke, Jen Henderson, Chris Hrobak, Jennifer Neidig, Jamie Snavelly, Sean Woznicki
Cover Design: Lisa Tosheff
Printed at: Yurchak Printing Inc.

Published in the United States of America by
Information Science Reference (an imprint of IGI Global)
701 E. Chocolate Avenue, Suite 200
Hershey PA 17033
Tel: 717-533-8845
Fax: 717-533-8661
E-mail: cust@igi-global.com
Web site: <http://www.igi-global.com/reference>

and in the United Kingdom by
Information Science Reference (an imprint of IGI Global)
3 Henrietta Street
Covent Garden
London WC2E 8LU
Tel: 44 20 7240 0856
Fax: 44 20 7379 0609
Web site: <http://www.eurospanbookstore.com>

Copyright © 2009 by IGI Global. All rights reserved. No part of this publication may be reproduced, stored or distributed in any form or by any means, electronic or mechanical, including photocopying, without written permission from the publisher.

Product or company names used in this set are for identification purposes only. Inclusion of the names of the products or companies does not indicate a claim of ownership by IGI Global of the trademark or registered trademark.

Library of Congress Cataloging-in-Publication Data

Encyclopedia of data warehousing and mining / John Wang, editor. -- 2nd ed.
p. cm.

Includes bibliographical references and index.

Summary: "This set offers thorough examination of the issues of importance in the rapidly changing field of data warehousing and mining"--Provided by publisher.

ISBN 978-1-60566-010-3 (hardcover) -- ISBN 978-1-60566-011-0 (ebook)

1. Data mining. 2. Data warehousing. I. Wang, John,

QA76.9.D37E52 2008

005.74--dc22

2008030801

British Cataloguing in Publication Data

A Cataloguing in Publication record for this book is available from the British Library.

All work contributed to this encyclopedia set is new, previously-unpublished material. The views expressed in this encyclopedia set are those of the authors, but not necessarily of the publisher.

If a library purchased a print copy of this publication, please go to <http://www.igi-global.com/agreement> for information on activating the library's complimentary electronic access to this publication.

Intelligent Query Answering

Zbigniew W. Ras

University of North Carolina, Charlotte, USA

Agnieszka Dardzinska

Bialystok Technical University, Poland

INTRODUCTION

One way to make Query Answering System (QAS) intelligent is to assume a hierarchical structure of its attributes. Such systems have been investigated by (Cuppens & Demolombe, 1988), (Gal & Minker, 1988), (Gaasterland et al., 1992) and they are called cooperative. Any attribute value listed in a query, submitted to cooperative QAS, is seen as a node of the tree representing that attribute. If QAS retrieves no objects supporting query q , from a queried information system S , then any attribute value listed in q can be generalized and the same the number of objects supporting q in S can increase. In cooperative systems, these generalizations are controlled either by users (Gal & Minker, 1988), or by knowledge discovery techniques (Muslea, 2004).

If QAS for S collaborates and exchanges knowledge with other systems, then it is also called intelligent. In papers (Ras & Dardzinska, 2004, 2006), a guided process of rules extraction and their goal-oriented exchange among systems is proposed. These rules define foreign attribute values for S and they are used to construct new attributes and/or impute null or hidden values of attributes in S . By enlarging the set of attributes from which queries for S can be built and by reducing the incompleteness of S , we not only enlarge the set of queries which QAS can successfully handle but also we increase the overall number of retrieved objects.

So, QAS based on knowledge discovery has two classical scenarios which need to be considered:

- **System is standalone and incomplete.**

Classification rules are extracted and used to predict what values should replace null values before any query is answered.

- **System is distributed with autonomous sites (including site S). User needs to retrieve objects from S satisfying query q containing nonlocal attributes for S .**

We search for definitions of these non-local attributes at remote sites for S and use them to approximate q (Ras & Zytkow, 2000), (Ras & Dardzinska, 2004, 2006).

The goal of this article is to provide foundations and basic results for knowledge-discovery based QAS.

BACKGROUND

Modern query answering systems area of research is related to enhancements of query-answering systems into intelligent systems. The emphasis is on problems in users posing queries and systems producing answers. This becomes more and more relevant as the amount of information available from local or distributed information sources increases. We need systems not only easy to use but also intelligent in handling the users' needs. A query-answering system often replaces human with expertise in the domain of interest, thus it is important, from the user's point of view, to compare the system and the human expert as alternative means for accessing information.

A knowledge system is defined as an information system S coupled with a knowledge base KB which is simplified in (Ras & Zytkow, 2000), (Ras & Dardzinska, 2004, 2006) to a set of rules treated as definitions of attribute values. If information system is distributed with autonomous sites, these rules can be extracted either locally from S (query was submitted to S) or from its remote sites. The initial alphabet of QAS associated with S contains all values of attributes in S , called local, and all decision values used in rules from KB . When KB is updated (new rules are added or some deleted), the alphabet for the local query answering

system is automatically changed. It is often assumed that knowledge bases for all sites are initially empty. Collaborative information system (Ras & Dardzinska, 2004, 2006) learns rules describing values of incomplete attributes and attributes classified as foreign for its site called a client. These rules can be extracted at any site but their condition part should use, if possible, only terms which can be processed by the query answering system associated with the client. When the time progresses more and more rules can be added to the local knowledge base which means that some attribute values (decision parts of rules) foreign for the client are also added to its local alphabet. The choice of which site should be contacted first, in search for definitions of foreign attribute values, is mainly based on the number of attribute values common for the client and server sites. The solution to this problem is given in (Ras & Dardzinska, 2006).

MAIN THRUST

The technology dimension will be explored to help clarify the meaning of intelligent query answering based on knowledge discovery and chase.

Intelligent Query Answering for Standalone Information System

QAS for an information system is concerned with identifying all objects in the system satisfying a given description. For example an information system might contain information about students in a class and classify them using four attributes of “hair color”, “eye color”, “gender” and “size”. A simple query might be to find all students with brown hair and blue eyes. When information system is incomplete, students having brown hair and unknown eye color can be handled by either including or excluding them from the answer to the query. In the first case we talk about optimistic approach to query evaluation while in the second case we talk about pessimistic approach. Another option to handle such a query would be to discover rules for eye color in terms of the attributes hair color, gender, and size. These rules could then be applied to students with unknown eye color to generate values that could be used in answering the query. Consider that in our example one of the generated rules said:

$(\text{hair, brown}) \wedge (\text{size, medium}) \rightarrow (\text{eye, brown})$.

Thus, if one of the students having brown hair and medium size has no value for eye color, then the query answering system should not include this student in the list of students with brown hair and blue eyes. Attributes hair color and size are classification attributes and eye color is the decision attribute.

We are also interested in how to use this strategy to build intelligent QAS for incomplete information systems. If query is submitted to information system S, the first step of QAS is to make S as complete as possible. The approach proposed in (Dardzinska & Ras, 2005) is to use not only functional dependencies to chase S (Atzeni & DeAntonellis, 1992) but also use rules discovered from a complete subsystem of S to do the chasing.

In the first step, intelligent QAS identifies all incomplete attributes used in a query. An attribute is incomplete in S if there is an object in S with incomplete information on this attribute. The values of all incomplete attributes are treated as concepts to be learned (in a form of rules) from S.

Incomplete information in S is replaced by new data provided by Chase algorithm based on these rules. When the process of removing incomplete values in the local information system is completed, QAS finds the answer to query in a usual way.

Intelligent Query Answering for Distributed Autonomous Information Systems

Semantic inconsistencies are due to different interpretations of attributes and their values among sites (for instance one site can interpret the concept “young” differently than other sites). Different interpretations are also due to the way each site is handling null values. Null value replacement by values suggested either by statistical or knowledge discovery methods is quite common before user query is processed by QAS.

Ontology (Guarino, 1998), (Van Heijst et al., 1997) is a set of terms of a particular information domain and the relationships among them. Currently, there is a great deal of interest in the development of ontologies to facilitate knowledge sharing among information systems.

Ontologies and inter-ontology relationships between them are created by experts in corresponding domain,

but they can also represent a particular point of view of the global information system by describing customized domains. To allow intelligent query processing, it is often assumed that an information system is coupled with some ontology. Inter-ontology relationships can be seen as semantical bridges between ontologies built for each of the autonomous information systems so they can collaborate and understand each other.

In (Ras and Dardzinska, 2004), the notion of optimal rough semantics and the method of its construction have been proposed. Rough semantics can be used to model semantic inconsistencies among sites due to different interpretations of incomplete values of attributes. Distributed chase (Ras and Dardzinska, 2006) is a chase-type algorithm, driven by a client site of a distributed information system DIS, which is similar to chase algorithms based on knowledge discovery and presented in (Dardzinska and Ras, 2005). Distributed chase has one extra feature in comparison to other chase-type algorithms: the dynamic creation of knowledge bases at all sites of DIS involved in the process of solving a query submitted to the client site of DIS.

The knowledge base at the client site may contain rules extracted from the client information system and also rules extracted from information systems at remote sites in DIS. These rules are dynamically updated through the incomplete values replacement process (Ras and Dardzinska, 2004, 2006).

Although the names of attributes are often the same among sites, their semantics and granularity levels may differ from site to site. As the result of these differences, the knowledge bases at the client site and at remote sites have to satisfy certain properties in order to be applicable in a distributed chase.

So, assume that system $S = (X, A, V)$, which is a part of DIS, is queried by user.

Chase algorithm, to be applicable to S , has to be based on rules from the knowledge base D associated with S which satisfies the following conditions:

1. Attribute value used in decision part of a rule from D has the granularity level either equal to or finer than the granularity level of the corresponding attribute in S .
2. The granularity level of any attribute used in the classification part of a rule from D is either equal or softer than the granularity level of the corresponding attribute in S .

3. Attribute used in the decision part of a rule from D either does not belong to A or is incomplete in S .

Assume again that $S = (X, A, V)$ is an information system (Pawlak, 1991), where X is a set of objects, A is a set of attributes (seen as partial functions from X into $2^{(V \times [0,1])}$ and, V is a set of values of attributes from A . By $[0,1]$ we mean the set of real numbers from 0 to 1. Let $L(D) = \{[t \rightarrow v_c] \in D: c \in \text{In}(A)\}$ be a set of all rules (called a knowledge-base) extracted initially from the information system S by ERID (Dardzinska and Ras, 2006), where $\text{In}(A)$ is a set of incomplete attributes in S .

Assume now that query $q(B)$ is submitted to system $S = (X, A, V)$, where B is the set of all attributes used in $q(B)$ and that $A \cap B \neq \emptyset$. All attributes in $B - [A \cap B]$ are called foreign for S . If S is a part of a distributed information system, definitions of foreign attributes for S can be extracted at its remote sites. Clearly, all semantic inconsistencies and differences in granularity of attribute values among sites have to be resolved first. In (Ras and Dardzinska, 2004) only different granularity of attribute values and different semantics related to different interpretations of incomplete attribute values among sites have been considered.

In (Ras and Dardzinska, 2006), it was shown that query $q(B)$ can be processed at site S by discovering definitions of values of attributes from $B - [A \cap B]$ at the remote sites for S and next use them to answer $q(B)$.

Foreign attributes for S in B , can be also seen as attributes entirely incomplete in S , which means values (either exact or partially incomplete) of such attributes should be ascribed by chase to all objects in S before query $q(B)$ is answered. The question remains, if values discovered by chase are really correct?

Classical approach, to this kind of problems, is to build a simple DIS environment (mainly to avoid difficulties related to different granularity and different semantics of attributes at different sites). As the testing data set we have taken 10,000 tuples randomly selected from a database of some insurance company in Charlotte, NC. This sample table, containing 100 attributes, was randomly partitioned into four subtables of equal size containing 2,500 tuples each. Next, from each of these subtables 40 attributes (columns) have been randomly removed leaving four data tables of the size 2,500×60 each. One of these tables was called a

client and the remaining 3 have been called servers. Now, for all objects at the client site, values of one of the attributes, which was chosen randomly, have been hidden. This attribute is denoted by d . At each server site, if attribute d was listed in its domain schema, descriptions of d using See5 software (data are complete so it was not necessary to use ERID) have been learned. All these descriptions, in the form of rules, have been stored in the knowledge base of the client. Distributed Chase was applied to predict what is the real value of the hidden attribute for each object x at the client site. The threshold value $\lambda = 0.125$ was used to rule out all values predicted by distributed Chase with confidence below that threshold. Almost all hidden values (2476 out of 2500) have been discovered correctly (assuming $\lambda = 0.125$) (Ras & Dardzinska, 2006).

Distributed Chase and Security Problem of Hidden Attributes

Assume now that an information system $S=(X,A,V)$ is a part of DIS and attribute $b \in A$ has to be hidden. For that purpose, we construct $S_b=(X,A,V)$ to replace S , where:

1. $a_s(x) = a_{S_b}(x)$, for any $a \in A - \{b\}$, $x \in X$,
2. $b_{S_b}(x)$ is undefined, for any $x \in X$,
3. $b_s(x) \in V_b$.

Users are allowed to submit queries to S_b and not to S . What about the information system $\text{Chase}(S_b)$? How it differs from S ?

If $b_s(x) = b_{\text{Chase}(S_b)}(x)$, where $x \in X$, then values of additional attributes for object x have to be hidden in S_b to guarantee that value $b_s(x)$ can not be reconstructed by Chase. Algorithm SCIKD for protection of sensitive data against Chase was proposed in (Im & Ras, 2007).

FUTURE TRENDS

One of the main problems related to semantics of an incomplete information system S is the freedom how new values are constructed to replace incomplete values in S , before any rule extraction process begins. This replacement of incomplete attribute values in some of the slots in S can be done either by chase or/and by a number of available statistical methods. This implies

that semantics of queries submitted to S and driven (defined) by query answering system QAS based on chase may often differ. Although rough semantics can be used by QAS to handle this problem, we still have to look for new alternate methods.

Assuming different semantics of attributes among sites in DIS, the use of global ontology or local ontologies built jointly with inter-ontology relationships among them seems to be necessary for solving queries in DIS using knowledge discovery and chase. Still a lot of research has to be done in this area.

CONCLUSION

Assume that the client site in DIS is represented by partially incomplete information system S . When a query is submitted to S , its query answering system QAS will replace S by $\text{Chase}(S)$ and next will solve the query using, for instance, the strategy proposed in (Ras and Dardzinska, 2004). Rules used by Chase can be extracted from S or from its remote sites in DIS assuming that all differences in semantics of attributes and differences in granularity levels of attributes are resolved first. We can argue here why the resulting information system obtained by Chase can not be stored aside and reused when a new query is submitted to S ? If system S is not frequently updated, we can do that by keeping a copy of $\text{Chase}(S)$ and next reusing that copy when a new query is submitted to S . But, the original information system S still has to be kept so when user wants to enter new data to S , they can be stored in the original system. System $\text{Chase}(S)$, if stored aside, can not be reused by QAS when the number of updates in the original S exceeds a given threshold value. It means that the new updated information system S has to be chased again before any query is answered by QAS.

REFERENCES

- Atzeni, P., DeAntonellis, V. (1992). *Relational Database Theory*, The Benjamin Cummings Publishing Company.
- Cuppens, F., Demolombe, R. (1988). Cooperative answering: a methodology to provide intelligent access to databases, *Proceedings of the Second International Conference on Expert Database Systems*, 333-353.

Dardzinska, A., Ras, Z. (2005). CHASE-2: Rule based chase algorithm for information systems of type lambda, *Post-proceedings of the Second International Workshop on Active Mining (AM'2003)*, Maebashi City, Japan, LNAI 3430, Springer, 258-270

Dardzinska, A., Ras, Z. (2006). Extracting rules from incomplete decision systems: System ERID, in *Foundations and Novel Approaches in Data Mining*, Studies in Computational Intelligence 9, Springer, 143-154

Gal, A., Minker, J. (1988). Informative and cooperative answers in databases using integrity constraints, *Natural Language Understanding and Logic Programming*, North Holland, 277-300.

Gaasterland, T., Godfrey, P., Minker, J. (1992). Relaxation as a platform for cooperative answering, *Journal of Intelligent Information Systems 1 (3)*, 293-321.

Giannotti, F., Manco, G. (2002). Integrating data mining with intelligent query answering,

in *Logics in Artificial Intelligence*, LNCS 2424, 517-520

Guarino, N., ed. (1998). *Formal ontology in information systems*, IOS Press, Amsterdam.

Im, S., Ras, Z. (2007). Protection of sensitive data based on reducts in a distributed knowledge discovery system, *Proceedings of the International Conference on Multimedia and Ubiquitous Engineering (MUE 2007)*, in Seoul, South Korea, IEEE Computer Society, 762-766

Muslea, I. (2004). Machine Learning for Online Query Relaxation, *Proceedings of KDD-2004*, in Seattle, Washington, ACM, 246-255

Pawlak, Z. (1991). *Rough sets-theoretical aspects of reasoning about data*, Kluwer.

Ras, Z., Dardzinska, A. (2004). Ontology based distributed autonomous knowledge systems, *Information Systems International Journal 29 (1)*, Elsevier, 47-58

Ras, Z., Dardzinska, A. (2006). Solving failing queries through cooperation and collaboration, *World Wide Web Journal 9(2)*, Springer, 173-186

Ras, Z., Zytkow, J.M. (2000). Mining for attribute definitions in a distributed two-layered DB system, *Journal of Intelligent Information Systems 14 (2/3)*, Kluwer, 115-130

Ras, Z., Zhang, X., Lewis, R. (2007). MIRAI: Multi-hierarchical, FS-tree based music information retrieval system, *Proceedings of RSEISP 2007*, LNAI 4585, Springer, 80-89

Van Heijst, G., Schreiber, A., Wielinga, B. (1997). Using explicit ontologies in KBS development, *International Journal of Human and Computer Studies 46, (2/3)*, 183-292.

KEY TERMS

Autonomous Information System: Information system existing as an independent entity.

Chase: Kind of a recursive strategy applied to a database V , based on functional dependencies or rules extracted from V , by which a null value or an incomplete value in V is replaced by a new more complete value.

Distributed Chase: Kind of a recursive strategy applied to a database V , based on functional dependencies or rules extracted both from V and other autonomous databases, by which a null value or an incomplete value in V is replaced by a new more complete value. Any differences in semantics among attributes in the involved databases have to be resolved first.

Intelligent Query Answering: Enhancements of query-answering systems into sort of intelligent systems (capable or being adapted or molded). Such systems should be able to interpret incorrectly posed questions and compose an answer not necessarily reflecting precisely what is directly referred to by the question, but rather reflecting what the intermediary understands to be the intention linked with the question.

Knowledge Base: A collection of rules defined as expressions written in predicate calculus. These rules have a form of associations between conjuncts of values of attributes.

Ontology: An explicit formal specification of how to represent objects, concepts and other entities that are assumed to exist in some area of interest and relationships holding among them. Systems that share the same ontology are able to communicate about domain of discourse without necessarily operating on a globally shared theory. System commits to ontology if its

observable actions are consistent with the definitions in the ontology.

Query Semantics: The meaning of a query with an information system as its domain of interpretation. Application of knowledge discovery and Chase in query evaluation makes semantics operational.

Semantics: The meaning of expressions written in some language, as opposed to their syntax which describes how symbols may be combined independently of their meaning.