

Optimization Algorithm and Data Security Problem in Distributed Information Systems

Agnieszka Dardzinska-Glebocka
Bialystok Technical University
Department of Computer Science
Wiejska 45A, 15351 Bialystok Poland
agnieszka.dardzinska@gmail.com

Abstract

In this work we present a new rule-discovery method for Distributed Information System. DIS is the system that connects information systems using network communication technology. This communication can be driven by request for knowledge needed to predict for maximal optimization which missing values can be replaced. In this work we recall the notion of a distributed information system to talk about handling semantic inconsistencies between sites. Semantic inconsistencies are due to different interpretations of attributes and their values on the concepts level among sites. Different interpretations can be also linked with a different way of treating null values among sites. Some attributes might be just hidden because of the security reason. In such case we have to be certain that the missing data can not be reconstructed from the available data by any known null value imputation method and that some information in IS can not be uncovered as well. Assuming that one attribute is hidden at one of the sites of DIS we will try to reconstruct this attribute. In this paper we will also show which values have to be hidden from users to guarantee that the hidden attribute can not be reconstructed.

1. Introduction

Distributed Information System (DIS) is a system that connects a number of information systems using network communication technology. In this paper, we assume that these systems are autonomous and incomplete.

Definition 1

By an Information System we mean a triple $S = (X, A, V)$ where:

X is a nonempty, finite set of objects;

A is a nonempty, finite set of attributes;

$V = \cup\{V_a : a \in A\}$ is a set of values of attributes, where

V_a is a set of values of attribute a , for any $a \in A$.

Additionally we assume that:

two sets of attribute values are disjoint $V_a \cap V_b = \emptyset$ for any two different attributes $a, b \in A$, and

$a : X \rightarrow V_a$ is a function for every $a \in A$.

We assume also, that the sum of the weights assigned to the attribute values has to be equal 1 in one tuple. In a case when there is an empty space in one tuple, it is interpreted in such way, that there can be all the values of attribute (from the domain of a given attribute) with equal weights. The definition of an information system of type λ and distributed information system (DIS) used in this paper was given in [7]. The type λ is used to check the weights assigned to values of attributes by Chase algorithm [7], if they are greater than or equal to given threshold λ . If the weight assigned by Chase to one of the attribute values is less than the given value λ , then this attribute value automatically is ruled out. Semantic inconsistencies among sites are due to different interpretations of attributes and their values among sites (for instance one site can interpret the concept e.g. beautiful, young or tall, in different way than another one). Ontology [1], [3], [4], [5], [9], [10], [11], [12] is understood as a set of terms of a particular information domain and the relationships between them. If two information systems agree on the ontology associated with attribute e.g. beautiful and its values, then such attribute can be used as a kind of semantical bridge between these systems. Different interpretations are also due to the way each site is handling null values. Null value replacement by a value predicted by statistical or some rule-based methods is quite common before queries are answered by query answering system. In [7], the notion of rough semantics and a method of its construction was proposed. The rough semantics can be used to optimize the model and handle semantic inconsistencies among sites due to different interpretations of incomplete values.

2. Query Processing with Incomplete Data

Let us start with the definition of partially incomplete information system S .

Definition 2

By a partially incomplete Information System $S = (X, A, V)$ of type λ we mean the incomplete information system, with three conditions:

- $(\forall x \in X)(\forall A \in A)a_S(x)$ is defined
- $[(a_S(x) = \{(a_i, p_i) : 1 \leq i \leq m\}) \rightarrow \sum_i p_i = 1]$
- $[(a_S(x) = \{(a_i, p_i) : 1 \leq i \leq m\}) \rightarrow (\forall i)(p_i) \geq \lambda]$

Now, let us assume that S_1, S_2 are partially incomplete information systems, both of type λ . The same objects from the set of objects X are stored in both systems and the same attributes from the set of attributes A are used to describe them. The meaning and granularity of values of attributes from A in both systems S_1, S_2 is also the same.

Additionally, we assume that

$$a_{S_1}(x) = \{(a_{1i}, p_{1i}) : 1 \leq p_i\} \text{ and}$$

$$a_{S_2}(x) = \{(a_{2i}, p_{2i}) : 1 \leq p_i\}.$$

We say that containment relation Ψ holds between S_1 and S_2 , if the following two conditions hold:

- $(\forall x \in X)(\forall a \in A)[card(a_{S_1}(x)) \geq card(a_{S_2}(x))],$
- $(\forall x \in X)(\forall a \in A)[card(a_{S_1}(x)) = card(a_{S_2}(x))] \rightarrow$

$$\rightarrow \sum_{i \neq j} |p_{2i} - p_{2j}| \succ \sum_{i \neq j} |p_{1i} - p_{1j}|.$$

If containment relation Ψ holds between S_1 and S_2 both of type λ , we also say that information system S_1 was mapped onto S_2 by containment mapping Ψ and denote that fact as

$$\Psi(S_1) = S_2$$

which means that

$$(\forall x \in X)(\forall a \in A)[\Psi(a_{S_1}(x)) = \Psi(a_{S_2}(x))].$$

We also say that containment relation Ψ holds between $a_{S_1}(x)$ and $a_{S_2}(x)$, for any $x \in X$, and $a \in A$.

If containment mapping Ψ converts an information system S_1 to S_2 , then we say that S_2 is more complete than S_1 . It means, that for a minimum one pair $(a, x) \in A \times X$, either Ψ has to decrease the number of attribute values in $a_S(x)$ or the average difference between confidences assigned to attribute values in $a_S(x)$ has to be increased by Ψ .

Example

Let us take two information systems S_1, S_2 both of type λ , represented as Figure 1 and Figure 2.

Figure 1. Incomplete information system S_1 .

X	a	b	c	d	e
x_1	$(a_1, \frac{1}{3})$ $(a_2, \frac{2}{3})$	$(b_1, \frac{2}{3})$ $(b_2, \frac{1}{3})$	c_1	d_1	$(e_1, \frac{1}{2})$ $(e_2, \frac{1}{2})$
x_2	$(a_2, \frac{1}{4})$ $(a_3, \frac{3}{4})$	$(b_1, \frac{1}{3})$ $(b_2, \frac{2}{3})$		d_2	e_1
x_3		b_2	$(c_1, \frac{1}{2})$ $(c_3, \frac{1}{2})$	d_2	e_3
x_4	a_3		c_2	d_1	$(e_1, \frac{2}{3})$ $(e_2, \frac{1}{3})$
x_5	$(a_1, \frac{2}{3})$ $(a_2, \frac{1}{3})$	b_1	c_2		e_1
x_6	a_2	b_2	c_3	d_2	$(e_2, \frac{1}{3})$ $(e_3, \frac{2}{3})$
x_7	a_2	$(b_1, \frac{1}{4})$ $(b_2, \frac{3}{4})$	$(c_1, \frac{1}{3})$ $(c_2, \frac{2}{3})$	d_2	e_2
x_8		b_2	c_1	d_1	e_3

Figure 2. Incomplete information system S_2 .

X	a	b	c	d	e
x_1	$(a_1, \frac{1}{3})$ $(a_2, \frac{2}{3})$	$(b_1, \frac{2}{3})$ $(b_2, \frac{1}{3})$	c_1	d_1	$(e_1, \frac{1}{2})$ $(e_2, \frac{1}{2})$
x_2	$(a_2, \frac{1}{4})$ $(a_3, \frac{3}{4})$	b_2	$(c_1, \frac{1}{3})$ $(c_2, \frac{2}{3})$	d_2	e_1
x_3	a_1	b_2	$(c_1, \frac{1}{2})$ $(c_3, \frac{1}{2})$	d_2	e_3
x_4	a_3		c_2	d_1	e_2
x_5	$(a_1, \frac{2}{3})$ $((a_2, \frac{1}{3})$	b_1	c_2		e_1
x_6	a_2	b_2	c_3	d_2	$(e_2, \frac{1}{3})$ $(e_3, \frac{2}{3})$
x_7	a_2	$(b_1, \frac{1}{4})$ $(b_2, \frac{3}{4})$	c_1	d_2	e_2
x_8	$(a_1, \frac{2}{3})$ $(a_2, \frac{1}{3})$	b_2	c_1	d_1	e_3

For explanation the Definition 2, let us look at the values of attribute a in both systems.

Assume that for a given object x :

$$a_{S_1}(x) = \{(a_1, \frac{1}{3}), (a_2, \frac{2}{3})\} \text{ and}$$

$$a_{S_2}(x) = \{(a_1, \frac{1}{5}), (a_2, \frac{4}{5})\}.$$

Clearly S_2 is closer to a complete system than S_1 with respect to $a(x)$, since uncertainty in the value of attribute a for x is lower in S_2 than in S_1 . It means that the containment mapping Ψ converts system S_1 to S_2 .

3. Query Processing with Distributed Data and Chase

Assume now that the knowledge-base $L(D)$ described as:

$$L(D) = \{(t \rightarrow v_c) \in D : c \in In(A)\}$$

is the set of all rules extracted from S by $ERID(S, \lambda_1, \lambda_2)$, where $In(A)$ is the set of incomplete attributes in S and there are given two thresholds describing minimum support and minimum confidence.

ERID is the algorithm for discovering rules from incomplete information systems which was presented in [2]. The type of incompleteness in [2] is the same as in this paper.

Assume now that a query $q(B)$ is submitted to system $S = (X, A, V)$, where B is the set of all attributes used in $q(B)$ and $A \cap B \neq \emptyset$. Attributes belonging to the set $B \setminus (A \cap B)$ are called hidden attributes in information system S . Hidden attributes for S can be seen as attributes which are entirely incomplete in system, which means exact or partially incomplete values of such attributes have to be ascribed to all objects in S . Stronger the consensus among sites on a value to be ascribed to x , better the result of the ascription process for x can be expected. Assuming that systems S_1, S_2 store the same sets of objects and use the same attributes to describe them, system S_1 is finer than system S_2 , if $\Psi(S_2) = S_1$.

Let us assume that $S = (X, A, V)$ is an information system of type λ and t is a term constructed in a standard way from values of attributes in V seen as constants and from two functors $+$ and $*$.

By $N_S(t)$ we mean the standard interpretation of a term t in S defined as in [6].

$N_S(t)$ is the standard interpretation of term t in $S = (X, A, V)$, and:

- $N_S(v) = \{(x, p) : (v, p) \in a(x)\}$, for any $v \in V_a$,
- $N_S(t_1 + t_2) = N_S(t_1) \oplus N_S(t_2)$,
- $N_S(t_1 * t_2) = N_S(t_1) \otimes N_S(t_2)$,

where for any

$$N_S(t_1) = \{(x_i, p_i)\}_{i \in I}, N_S(t_2) = \{(x_j, q_j)\}_{j \in J} \text{ we have:}$$

- $N_S(t_1) \oplus N_S(t_2) = \{(x_j, p_j)\}_{j \in J \setminus I} \cup \{(x_i, p_i)\}_{i \in I \setminus J} \cup \{(x_i, \max(p_i, q_i))\}_{i \in I \cap J}$
- $N_S(t_1) \otimes N_S(t_2) = \{(x_i, p_i \cdot q_i)\}_{i \in I \cap J}$

The incomplete value imputation algorithm Chase [8], based on the above semantics converts information system S of type λ to a more complete, new information system of the same type. Algorithm ERID can be used to extract rules from the first information system and next can be applied in Chase.

4. Security Problem of Hidden Attributes

Assume that information system S is a distributed information system, and the attribute $h \in A$ is hidden. We also assume, that $S_h = (X, A, V)$, where

- $(\forall a \in A \setminus \{h\})(\forall x \in X) a_S(x) = a_{S_h}(x)$
- $(\forall x \in X) h_{S_h}$ is undefined
- $h_{S_h} \in V_h$

The assumption is, that the user can only submit a query to S_h and not to S . We show how to find a minimal number of additional values which should be hidden to be sure that the values of attribute h can not be reconstructed by Chase for any $x \in X$.

Example

Let us take IS from Figure 1. Let this system be a system of type $\lambda = \frac{1}{4}$. Let us assume that attribute d is hidden in S .

Figure 3. Information system S_d .

X	a	b	c	d	e
x_1	$(a_1, \frac{1}{3})$ $(a_2, \frac{2}{3})$	$(b_1, \frac{2}{3})$ $(b_2, \frac{1}{3})$	c_1		$(e_1, \frac{1}{2})$ $(e_2, \frac{1}{2})$
x_2	$(a_2, \frac{1}{4})$ $(a_3, \frac{3}{4})$	$(b_1, \frac{1}{3})$ $(b_2, \frac{2}{3})$			e_1
x_3		b_2	$(c_1, \frac{1}{2})$ $(c_3, \frac{1}{2})$		e_3
x_4	a_3		c_2		$(e_1, \frac{2}{3})$ $(e_2, \frac{1}{3})$
x_5	$(a_1, \frac{2}{3})$ $(a_2, \frac{1}{3})$	b_1	c_2		e_1
x_6	a_2	b_2	c_3		$(e_2, \frac{1}{3})$ $(e_3, \frac{2}{3})$
x_7	a_2	$(b_1, \frac{1}{4})$ $(b_2, \frac{3}{4})$	$(c_1, \frac{1}{3})$ $(c_2, \frac{2}{3})$		e_2
x_8		b_2	c_1		e_3

Assume, the following rules were extracted at the remote sites for S_d :

- $r_1 : a_2 \cdot b_2 \rightarrow d_2$ with sup=3 and conf=1
- $r_2 : b_2 \cdot c_1 \rightarrow d_1$ with sup=2 and conf=1
- $r_3 : b_2 \cdot c_3 \rightarrow d_2$ with sup=3 and conf=1
- $r_4 : a_1 \cdot c_1 \rightarrow d_2$ with sup=3 and conf=1
- $r_5 : a_1 \cdot b_2 \rightarrow d_1$ with sup=2 and conf=1
- $r_6 : a_2 \cdot c_1 \rightarrow d_2$ with sup=4 and conf=1

All the above rules have confidence equal 1. Additionally rules r_2 and r_5 have support 2, rules r_1, r_3, r_4 have

support 3, and rule r_6 has support equals 4.

Let us consider the first tuple (for x_1) from S_d (Figure 3).

It supports rules r_1, r_2, r_4, r_5, r_6 .

Rule r_1 supports d_2 with weight $\frac{2}{3} \cdot \frac{2}{3} \cdot \frac{1}{3} \cdot 3 \cdot 1 = \frac{2}{3}$.

Rule r_2 supports d_1 with weight $\frac{2}{3} \cdot 1 \cdot 2 \cdot 1 = \frac{2}{3}$.

Rule r_4 supports d_2 with weight $\frac{2}{3} \cdot 1 \cdot 3 \cdot 1 = 1$.

Rule r_5 supports d_1 with weight $\frac{2}{3} \cdot \frac{1}{3} \cdot 2 \cdot 1 = \frac{2}{9}$.

Rule r_6 supports d_2 with weight $\frac{2}{3} \cdot 1 \cdot 4 \cdot 1 = \frac{2}{3}$.

We can calculate the total weight for d_1 , which is $\frac{2}{3} + \frac{2}{9} = \frac{8}{9}$, and the total weight for d_2 , which is $\frac{2}{3} + 1 + \frac{2}{3} = \frac{13}{3}$.

Because $\frac{8}{9} : \frac{13}{3} < \frac{1}{4}$, the value d_1 is rule out and the same can not be predicted by Chase. Therefore the values of attributes in this tuple do not have to be hidden

Following the similar strategy for all the objects from information system we obtain a new information system S_d (Figure 4)

Figure 4. New information system S_d .

X	a	b	c	d	e
x_1	$(a_1, \frac{1}{3})$ $(a_2, \frac{2}{3})$	$(b_1, \frac{2}{3})$ $(b_2, \frac{1}{3})$	c_1		$(e_1, \frac{1}{2})$ $(e_2, \frac{1}{2})$
x_2	$(a_2, \frac{1}{4})$ $(a_3, \frac{2}{4})$	$(b_1, \frac{1}{3})$ $(b_2, \frac{2}{3})$			e_1
x_3		b_2	$(c_1, \frac{1}{2})$ $(c_3, \frac{1}{2})$		e_3
x_4	a_3		c_2		$(e_1, \frac{2}{3})$ $(e_2, \frac{2}{3})$
x_5	$(a_1, \frac{2}{3})$ $((a_2, \frac{1}{3}))$	b_1	c_2		e_1
x_6					$(e_2, \frac{1}{3})$ $(e_3, \frac{2}{3})$
x_7		$(b_1, \frac{1}{4})$ $(b_2, \frac{3}{4})$	$(c_1, \frac{1}{3})$ $(c_2, \frac{2}{3})$		e_2
x_8		b_2	c_1		e_3

We can notice, that all the values of attributes a, b, c for tuple x_6 have been removed, as they are so weak, they can be helpful in prediction process for attribute value d . Then, the hidden attribute can not be reconstructed by Chase from the available data in S_d , for any object x .

Assume, that the knowledge base contains rules extracted in DIS at server sites for S_d with a goal to reconstruct hidden attribute d . For each object x , first of all we discover all rules supported by the tuple. We have to take into consideration all the possibilities, which are as follows:

- there is only one rule supported by object x in S_d
- there is a set of rules supported by object x in S_d .

In the first case, when we have one rule $r : t \rightarrow f$ supported

by object x , and $d = f$, the value f is predicted correctly by r , which means, that minimum one of the attributes listed in t has to be additionally hidden for x .

In the second case, there is a situation, where $r_1 : t_1 \rightarrow f_1, r_2 : t_2 \rightarrow f_1, \dots, r_i : t_i \rightarrow f_1$ is the set of all rules supported by x and $d = f$. In this case a minimal set of attributes covering all terms t_1, t_2, \dots, t_i needs to be additionally hidden for x in S .

There exists also the third possibility, where there is a set of rules, such that $r_1 : t_1 \rightarrow f_1, r_2 : t_2 \rightarrow f_2, \dots, r_i : t_i \rightarrow f_i$ supported by x . In this case we calculate support of the rule s_i , and its confidence c_i . Let λ be also given threshold for minimal confidence in attribute values for objects in S .

The confidence in attribute value e for x in S can be defined as follows:

$$ConfS(e, x) = \sum \{s_i \cdot c_i : 1 \leq i \leq k \wedge e = d_i / s_i \cdot c_i : 1 \leq i \leq k\}$$

In such case:

If $ConfS(f, x) > \lambda$ and $(\exists e \neq f)[ConfS(e, x) > \lambda]$, we do not have to hide any additional slots for x .

If $ConfS(f, x) > \lambda$ and $(\exists e \neq f)[ConfS(e, x) < \lambda]$, we have to hide additional slots for x .

If $ConfS(f, x) < \lambda$ and $(\exists e \neq f)[ConfS(e, x) > \lambda]$, we do not have to hide any additional slots for x .

5. Testing and Implementation

Query Answering System (QAS) for an incomplete information system was implemented on *Sparc20* using *ANSI C* and the *Rough Set Library*, ver. 2.0.

The primary purpose of this system was the testing of rule-based chase approach to answering queries for incomplete information systems. As the testing information system, *Zoo Database* from *UCIrvine database repository* has been chosen. This complete database containing 101 names of animals classifies them with respect to 17 Boolean-valued attributes and one type attribute. In order to test the effectiveness of our method on *Zoo database*, different percentage of randomly chosen attribute values in the information system S , to be replaced with unknown values, have been supplied. The answer to the same query is recorded, against the initial information system (which is complete), the incomplete information system (obtained after random removal of attribute values from S), and finally against the information system obtained as the result of Chase. After all query attributes have had certain rules generated and unknown values chased, then the entire incomplete information system is loaded into the data structure. Next, we

update the unknown attribute values in the incomplete information system with the chased values saved in the data structure. The updated information system is then checked for unknown query attribute values. When the information system is complete, with respect to the query attributes, no more passes are needed. The other stopping condition for the iterations of the rule-based chasing of unknown query attribute values occurs when a pass fails to discover any new values. This stopping condition is known as reaching a fixed point. Finally, the query runs against the incomplete information system updated with the chased query attribute values. The results of the query are recorded and compared to the query answer on the complete and incomplete information systems.

Using our method for predicting what attribute value should replace an incomplete value has a clear advantage over any other method for predicting incomplete values, mainly because of the use of existing associations between values of attributes. To find these associations we can use either any association rule mining algorithm or any rule discovery algorithm like *LERS* or *Rosetta*.

6. Conclusions

Presented method seems to be very interesting and promising in hiding some values of attributes from data security point of view. We showed the possibility and importance of hiding the attributes in information systems. For any tuple x we are able to identify all the rules supported by that tuple. On the basis of these rules, we calculate the total support for each value of the hidden attribute. These total supports are used to calculate the confidence in each of these values. If the confidence is below the given threshold λ , then such value is rule out. We need minimum two weighted values remaining if the correct value is one of them. This can be achieved by replacing some values by null ones. The suggested strategy provides a way to identify a minimal number of additional slots in IS required to be hidden if one of the attributes in IS has to be hidden.

6.1. References

[1] Benjamins, V. R., Fensel, D., Prez, A. G., Knowledge management through ontologies, Proceedings of the 2nd International Conference on Practical Aspects of Knowledge Management (PAKM-98), Basel, Switzerland, 1998

[2] Dardzinska, A., Ras, Z.W., On Rules Discovery from Incomplete Information Systems, Proceedings of ICDM'03 Workshop on Foundations and New Directions of Data Mining, (Eds: T.Y. Lin, X. Hu, S. Ohsuga, C.Liau),

Melbourne, Florida, IEEE Computer Society, 2003, 31-35

[3] Fensel, D., Ontologies: a silver bullet for knowledge management and electronic commerce, Springer-Verlag, 1998

[4] Guarino, N., Formal Ontology in Information System, IOS Press, Amsterdam, 1998.

[5] Guarino, N., Giaretta, P, Ontologies and knowledge bases towards a terminological clarification, Towards Very Large Knowledge Bases: Knowledge Building and Knowledge Sharing, IOS Press, 1995.

[6] Ras, Z.W., Joshi, S. Query approximate answering system for an incomplete DKBS, Fundamenta Informaticae Journal, IOS Press, Vol. 30, No. 3/4, 1997, 313-324

[7] Ras, Z.W, Dardzinska, A, Ontology Based Distributed Autonomous Knowledge Systems, Information Systems International Journal, Elsevier, Vol.29, No.1, 2004, 47-58.

[8] Ras, Z.W., Dardzinska, A. Query answering based on collaboration and chase, Proceedings of FQAS 2004 Conference, Lyon, France, LNCS/LNAI, No. 3055, Springer-Verlag, 125-136

[9] Sowa, J.F., Ontology, metadata and semiotics, B.Ganter G.W.ineau, Conceptual Structures: Logical, Linguistic, and Computational Issues, LNAI, No.1867, Springer-Verlag, Berlin, 2000, pp.55-81.

[10] Sowa, J.F., Knowledge Representation: Logical, Philosophical and Computational Foundations, Brooks/ColePublishing Co., Pacific Grove, CA, 2000.

[11] Sowa, J.F., Ontological categories, L. Albertazzi,ed., Shapes of Forms: From Gestalt Psychology and Phenomenology to Ontology and Mathematics, Kluwer Academic Publishers, Dordrecht, 1999, pp. 307-340.

[12] Van Heijst, G., Schreiber, A., Wielinga, B. Using explicit ontologies in KBS development, International Journal of Human and Computer Studies, Vol. 46, No. 2/3, 183-292.

7. Acknowledgement

The work on the paper has been supported by W/WI/15/07.