

Deriving Conceptual Schema from XML Databases

by Oviliani Yuliana

Submission date: 23-Apr-2023 08:45PM (UTC+0700)

Submission ID: 2072702060

File name: Deriving_Conceptual_Schema_from_XML_Databases.doc (2.25M)

Word count: 3197

Character count: 16087

Deriving Conceptual Schema from XML Databases

Oviliani Yenty Yuliana¹, Suphamit Chittayasothorn²

¹Department of Informatics Engineering, Petra Christian University, Surabaya, Indonesia

²Faculty of Engineering, King Mongkut's Institute of Technology Ladkrabang, Thailand
ovi@peter.petra.ac.id, suphamit@kmitl.ac.th

Abstract

In this paper, two concepts from different research areas are addressed together, namely functional dependency (FD) and multidimensional association rule (MAR). FD is a class of integrity constraints that have gained fundamental importance in relational database design. MAR is a class of patterns which has been studied rigorously in data mining. We employ MAR to mine the interesting rules from XML Databases. The mined interesting rules are considered as candidate FDs whose all confidence itemsets are 100%. To prune the weak rules, we pay attention to support and correlation itemsets. The final strong rules are used to generate an Object-Role Model conceptual schema diagram.

1. Introduction

In recent years, XML [1] has emerged as the dominant standard for representing and exchanging data over the Internet. However, there are many existing XML Databases which have been implemented without schemas and consistency checking. As a result, XML Databases may be inconsistent, incomplete or difficult to maintain. Such systems are also normally poor documented. These problems can be overcome by having XML Schemas which can be created using conceptual modeling techniques.

XML Schemas design using the Object-Role Modeling (ORM) as its conceptual schema are conducted by [2,3,4]. In addition, [4] captured all of ORM constraints that are still not defined in [1] and used XQueries to detect invalid constraints in XML Databases. Furthermore, [5] proposed reengineering the existing XML Databases using a conceptual schema approach.

At present, more and more scholars conduct research on association rules mining and many

methods are based on the Apriori Algorithm that was proposed by [6]. Single Association Rule Mining from XML Data was conducted by [7]. Moreover, [8] studied extracting Association Rules from XML Documents using XQuery [9]

In this paper, we employ Multidimensional Association Rule (MAR) to improve the ORM reverse engineering which was proposed by [5]. So far, deriving conceptual schemas from XML Databases using MAR has not been addressed. MAR is used for mining the interesting rules, i.e. candidate FDs, which all confidence itemsets are 100%. Usually FDs are based on superkey [10,11]. XQueries are applied for calculating confidence, support, and correlation itemsets. To prune the weak rules support and correlation itemsets are considered.

2. Basic concepts and notation

2.1. Functional dependencies

Let R be a relation schema. A subset K of R is a superkey of R if, in any legal relation $r(R)$, for all pairs t_1 and t_2 of tuples in r such that $t_1 \neq t_2$, then $t_1[K] \neq t_2[K]$. That is, no two tuples in any legal relation $r(R)$ may have the same value on attribute set K .

The notation of functional dependency generalizes the notion of superkey [10]. Let $\alpha \subseteq R$ and $\beta \subseteq R$. The functional dependency $\alpha \rightarrow \beta$ holds on R if, in any relation $r(R)$, for all pairs of tuples t_1 and t_2 in r such that $t_1[\alpha] = t_2[\alpha]$. It is also the case that $t_1[\beta] = t_2[\beta]$. Using the functional-dependency notation, K is a superkey of R if $K \rightarrow R$. That is, K is a superkey if, whenever $t_1[K] = t_2[K]$, it is also the case that $t_1[R] = t_2[R]$, that is $t_1 = t_2$. In this paper, we will mine FDs using the elementary facts concept which was proposed by [12].

2.2. Object-role modeling

ORM is called “fact-oriented modeling” because it expresses the information in terms of simple/elementary facts. An elementary fact is defined by [12] as an assertion that an object has a property, or that one or more objects participate in a relationship, where the fact cannot be split into simpler facts with the same object types without information loss.

In order to reverse engineering XML Databases into ORM conceptual schema, our concern is the second step of the conceptual schema design procedure, i.e. draw the fact types and apply a population check. In this paper, fact types are populated with fact instances from XML Databases using XQuery. For population checking, we propose an application that implements the MAR concept.

2.3. Multidimensional association rules

Let $J=\{i_1, i_2, \dots, i_m\}$ be a set of items. Let D , the task-relevant data, be a set of database transactions where each transaction T is a set of items such that $T \subseteq J$. Let A be a set of items. A transaction T is said to contain A if and only if $A \subseteq T$. [13] defined an association rule as an implication¹ of the form $A \Rightarrow B$, that is, $A_1 \wedge \dots \wedge A_m \Rightarrow B_1 \wedge \dots \wedge B_n$, where $A \subseteq J$, $B \subseteq J$, and $A \cap B = \emptyset$. Association rules that involve two or more dimensions or predicates can be referred to as Multidimensional Association Rules, for instance: $\text{age}(A, "20 \dots 29") \wedge \text{income}(A, "20K \dots 29K") \Rightarrow \text{buys}(X, "CD \text{ player}").$

The rule $A \Rightarrow B$ holds in the transaction set D with support s , where s is the percentage of transactions in D that contain $A \cup B$, i.e. both A and B . This is taken to be the probability, $P(A \cup B)$. The rule $A \Rightarrow B$ has confidence c in the transaction set D if c is the percentage of transactions in D containing A that also contain B . This is taken to be the conditional probability, $P(B/A)$.

$$\text{Support } (A \Rightarrow B) = P(A \cup B) = \frac{P(A \cup B)}{P(U)} \dots \dots \dots (1)$$

$$\text{Confidence } (A \Rightarrow B) = P(B/A) = \frac{P(A \cup B)}{P(A)} \dots \dots \dots (2)$$

Association Rules mined using a support-confidence framework are useful for many application. However, the support-confidence framework can be misleading in that it may identify a rule $A \Rightarrow B$ as interesting when the occurrence of A does not imply the occurrence of B . The occurrence of itemset A is independent of the occurrence of itemset B if $P(A \cup B) = P(A)P(B)$; otherwise itemsets A and B are dependent and correlated as events.

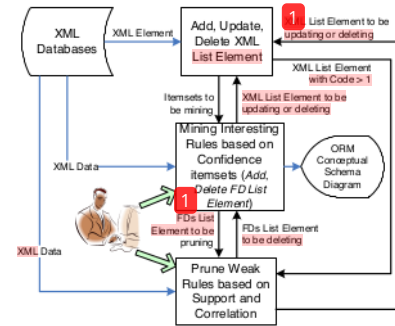
$$\text{Correlation } A, B = \frac{P(A \cup B)}{P(A)P(B)} \dots \dots \dots (3)$$

If the resulting value of (3) is less than 1, then the occurrence of A is negatively correlated with the occurrence of B . If the resulting value is greater than 1, then A and B are positively correlated, meaning the occurrence of one implies the occurrence of the other. If the resulting value is equal to 1, then A and B are independent and there is no correlation between them.

In this paper, a confidence value is used for mining the rule $A \Rightarrow B$ as the interesting FDs. If and only if all confidences in itemsets are 100% (see FD definition on section 2.1) then the rules $A \Rightarrow B$ is an interesting FDs. For FD purpose, we also mine the rule $B \Rightarrow A$. In addition, support and correlation are used to prune the interesting rules when more than one determinant determines the same dependant.

3. Mining multidimensional association rules from XML databases

A framework for deriving conceptual schema from XML Databases is shown in Figure 1. There are three main processes, i.e. modify XML List, mine interesting rules, and prune weak rules. The first process is to create XML List from XML Databases then generate itemsets to be mined by the second process. Create the rules (FDs) which are to be pruned in the third process and get inputs from the second and the third processes to update or delete the XML List elements.



¹ Figure1. The deriving conceptual schema from XML databases framework

The second process is to mine interesting rules by joining and calculating confidence itemsets using XQuery. Figure 2(a) and Figure 2(b) show XQueries for 2-itemsets and 3-itemsets respectively. The XQuery can be extended for further itemsets. The XQueries implement the confidence equation (2). The itemset is

an FD, if and only if all confidences of the itemset are 100% then save the itemset into an FD List.

The last process, prune weak rules (FDs) when XML List elements are referred more than one with considering average support and correlation itemsets.

Support equation (1) and correlation equation (3) are implemented on XQuery in Figure 2(c). Prune the lowest support FDs with the correlation FDs greater than or lower than one. The final FD List is used to generate an ORM conceptual schema diagram.

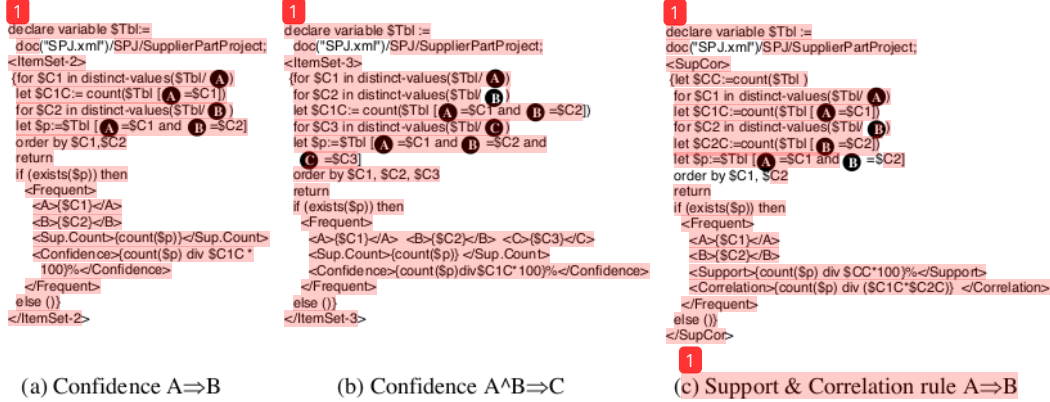


Figure 2. XQueries for the calculation of confidence, support, and correlation itemsets

Three struct data types in Figure 3, i.e. EXML, EDT, and EDP are proposed to support the framework. EXML is used in XML List for storing the third level XML elements in XML Databases. There are three possibilities value in $XML \rightarrow Code$, i.e. -1 (the element is referred to as determinant), 0 (the element is not referred), ≥ 1 (the element is referred at least one). EDT and EDP are used for constructing a FD List, i.e. DT List and DP List. DP List is used for storing dependent items. $DP \rightarrow next$ is used for linking between dependents. Set the last $DP \rightarrow next$ to Nil. DT List is used for storing determinant items that can be single or composite item. In case composite determinant, $DT \rightarrow nextT$ is used for linking between determinants. $DT \rightarrow nextB$ is used for linking DT elements with DP elements. In addition $DT \rightarrow top$ and $DT \rightarrow bottom$ is used for linking between DT elements. In case, no element in the top or bottom DT element than set $DT \rightarrow top$ or $DT \rightarrow bottom$ to Nil.

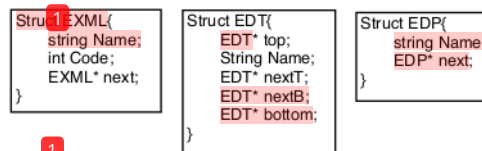


Figure 3. The proposed struct data type

An algorithm for mining and pruning the interesting FDs is as follows:

Step 1: Create an XML List. For every element on

the third level element in XML Databases, create an element in XML List. The XML List is a guide for generating itemsets, start from 2-itemsets, 3-itemsets, etc. until all $XML \rightarrow Code \neq 0$. For generating n-itemsets, n array of pointers are needed, namely $XML[i]$ for $i=0, \dots, n-1$. Set pointer $XML[0]$ to the first element XML List and set pointer $XML[1]$ to $XML[0] \rightarrow next$. Go to step 3.

Step 2: Move pointer $XML[0]$. Set pointer $PXML$ to $XML[0] \rightarrow next$. If $PXML \rightarrow next \neq Nil$ and still have $XML \rightarrow Code = 0$ then move pointer $XML[0]$ to $PXML$. Set pointer $XML[1]$ to $XML[0] \rightarrow next$. Otherwise go to step 8 for pruning the weak rules.

Step 3: Mine interesting rules 2-itemsets. Set false to Result1 and Result2. Pass $XML[0] \rightarrow Name$ and $XML[1] \rightarrow Name$ into the XQuery in Figure 2(a) for calculating confidence. If all confidences are 100% then set true to Result1. In addition, calculate confidence by passing $XML[1] \rightarrow Name$ and $XML[0] \rightarrow Name$ into the XQuery. If all confidences are 100% then set true to Result2.

Step 4: If both Result1 and Result2 are true, it means $XML[0]$ element is a determinant of $XML[1]$ element or the $XML[1]$ element is also a determinant of $XML[0]$. Ask the user to choose the right one. If the user chooses $XML[1]$ element is a determinant of $XML[0]$ element then go to step 7. Otherwise, check if $XML[0] \rightarrow Code = 0$ then create a DT element and fill $DT \rightarrow Name$ with $XML[0] \rightarrow Name$. Create a DP element and fill $DP \rightarrow Name$ with $XML[1] \rightarrow Name$. Moreover, set $XML[0] \rightarrow Code$ with -1. If

item in the itemsets, for instance $\{SNo, SCity\}$ and $\{SCity, SNo\}$. The confidence itemsets $\{SNo, SCity\}$ and $\{SCity, SNo\}$ are shown in Figure 6. All confidences for itemsets $\{SNo, SCity\}$ are 100% so assign true to Result1. It means rule $SNo \Rightarrow SCity$ is an interesting FD rule. However, all confidences for itemsets $\{SCity, SNo\}$ are not 100% so assign false to Result2. If only Result1 is true, then store item SNo into DT List and store item SCity into DP List. An example for only Result2 is true, when mining a rule

$PName \Rightarrow PNo$ and a rule $PNo \Rightarrow PName$. Only all confidences for itemsets $\{PNo, PName\}$ are 100%. Therefore, store PNo into DT List and PName into DP List. A ¹ example for both Result1 and Result2 are true, when mining rules $JName \Rightarrow JNo$ and $JNo \Rightarrow JName$. In this case, ask the user to choose only one rule. The XML List after processed 2-itemsets is shown in Figure 5(b) and The FD List is shown in Figure 8(a).

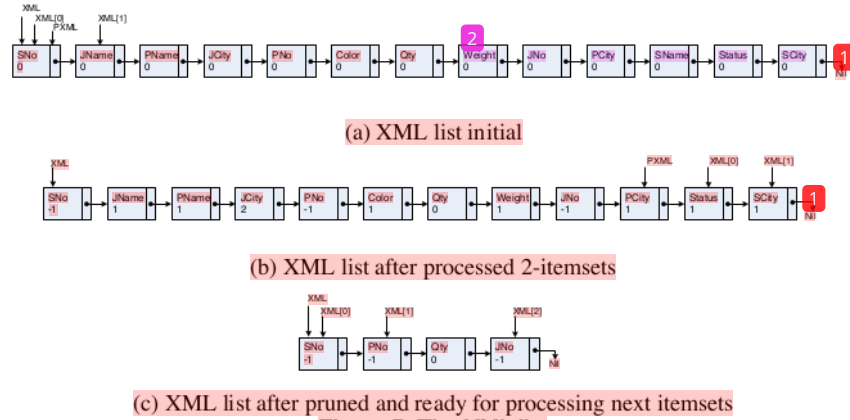


Figure 5. The XML list

```
<ItemSet-2>
<Frequent>
<A>S1</A>
<B>London</B>
<Sup.Count>2</Sup.Count>
<Confidence>100%</Confidence>
</Frequent>
... etc
</ItemSet-2>
```

SNo	SCity	Sup.count	Confidence
S1	London	2	100%
S2	Paris	8	100%
S3	Paris	2	100%
S4	London	2	100%
S5	Athens	10	100%

SCity	SNo	Sup.count	Confidence
Athens	S5	10	100%
London	S1	2	50%
London	S4	2	50%
Paris	S2	8	80%
Paris	S3	2	20%

Figure 6. The confidence itemsets $\{SNo, SCity\}$ and $\{SCity, SNo\}$

To prune the weak FD rules use XML List as a guide line, i.e. $XML \rightarrow Code \geq 1$. If $XML \rightarrow Code > 1$ then FD should be pruned according to the support and the correlation those are calculated by XQuery in Figure 2(c). For example, in Figure 5(b) element Name JCity, the Code is 2, i.e. itemsets $\{JName, JCity\}$ and $\{JNo, JCity\}$. The calculation support and correlation itemsets is show in Figure 7. The average support and the average correlation for rule $JName \Rightarrow JCity$ and

$JNo \Rightarrow JCity$ are same, i.e. 14% and 0.30 respectively. Therefore ask the user to choose the rule to be pruned. In this study case, we choose to prune the rule $JName \Rightarrow JCity$. Furthermore, delete all XML elements in XML List with Code are greater or equal one for preparing to next process mining. As a result the XML List is shown in Figure 5(c). There is still one XML List element with Code is 0.

```
<SupCor>
<Frequent>
<A>Console</A>
<B>Athens</B>
<Support>33.33%</Support>
<Correlation>0.1</Correlation>
</Frequent>
... etc
</SupCor>
```

JName	JCity	Support	Correlation
Console	Athens	33.33%	0.10
Display	Rome	20.83%	0.20
EDS	Oslo	8.33%	1.00
OCR	Athens	8.33%	0.10
RAID	London	8.33%	0.20
Sorter	Paris	12.50%	0.33
Tape	London	12.50%	0.20
Average		14.29%	0.30

JNo	JCity	Support	Correlation
J1	Paris	12.50%	0.33
J2	Rome	20.83%	0.20
J3	Athens	8.33%	0.10
J4	Athens	33.33%	0.10
J5	London	8.33%	0.20
J6	Oslo	4.17%	1.00
J7	London	12.50%	0.20
Average		14.29%	0.30

Figure 7. The support and correlation itemsets $\{JName, JCity\}$ and $\{JNo, JCity\}$

Mine 3-itemsets ($A \wedge B \Rightarrow C$) by combination two items A-B with Code -1 and one item C with Code at least 0, i.e. {SNo, PNo, Qty}, {SNo, JNo, Qty}, {PNo, JNo, Qty}. Pass 3-itemsets into the XQuery in Figure 2(B). However, no all confidence for every 3-itemsets is 100%. Therefore, mine 4-itemset ($A \wedge B \wedge C \Rightarrow D$), i.e. {SNo, PNo, JNo, Qty}. All confidences in

the itemsets are 100%. As a result store the rule in FD List and increment the Qty \rightarrow Code by one. If all XML \rightarrow Code is not 0 then stop the mining process. As a result FD List is shown Figure 8(b). The last step is used for generating an ORM conceptual schema in Figure 8(c), for detail algorithm refer to [5].

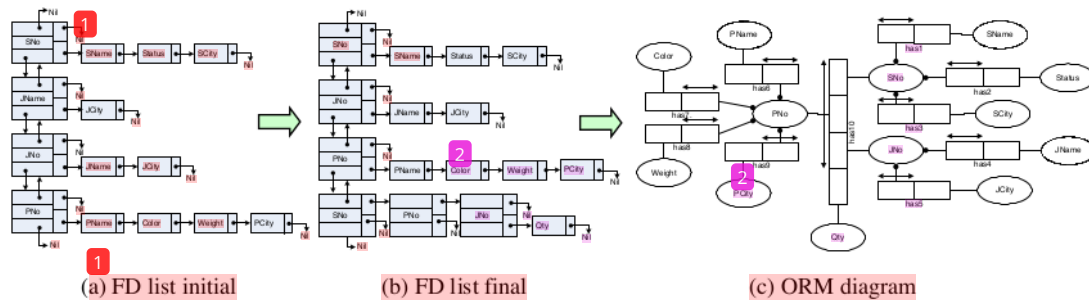


Figure 8. The FD list and the generated ORM diagram

5. Conclusion

This paper presents a conceptual schema reverse engineering technique and shows that MAR can be used for mining FDs in XML Databases. If XML Databases are not big enough, it is possible to produce interesting rules which depend on difference determinants. The support and the correlation itemsets can be employed to prune the weak FD rules. It is concluded that MAR techniques can improve our proposed conceptual schema reverse engineering technique [5] and can be used generally to reverse engineer FDs from XML Databases.

6. References

- [1] World Wide Web Consortium. XML 1.0 (Second Edition) W3C Recommendation. <http://www.w3.org/XML>.
- [2] Linda, Bird, Andrew Goodchild, and Terry Halpin. "Object Role Modeling and XML-Schema", *Proc. of ER 2000*, USA, October 2000, pp. 1-14.
- [3] Narudol C. and Suphamit C., "An Object and XML Database Schemas Design Tool", *Proc. of ITCC 2004*, USA, April 2004, pp. 421-424.
- [4] Yuliana, Oviliani Y. and Chittayasothorn S., "A Conceptual Schema Based XML Schema with Integrity Constraints", *Proc. of ICHIT 2008*, Korea, August 28-29 2008, pp. 19-24.
- [5] Yuliana, Oviliani Y. and Chittayasothorn S., "XML Schema Re-Engineering Using a Conceptual Schema Approach", *Proc. of ITCC 2005*, USA, April 4-6 2005, pp. 255-260.
- [6] R. Agrawal and Srikant R, "Fast Algorithm for Mining Association Rules", *Proc. of the 20th International Conference on Very Large Databases*, Chile, 1994, pp. 487-499.
- [7] Myint Myint Khaing and Nilar Thein, "An Efficient Association Rule Mining for XML Data", *Proc. of SICE-ISACE 2006*, Korea, October 18-21 2006, pp. 5782-5786.
- [8] Jacky W.W. Wan and Gillian Dobbie, "Extracting Association Rules from XML documents Using XQuery", *Proc of WIDM*, USA, November 7-8 2003, pp. 94-97.
- [9] D. Chamberlin, "XQuery: An XML query language", *IBM Systems Journal Vol. 41 No. 4*, 2002, pp. 597-615.
- [10] Date, C. J., *An Introduction to Database Systems the 7th edition*, Addison Wesley Longman, Inc., USA, 2000.
- [11] Pedro Sousa, Lurdes Pedro-de-Jesus et. al, "Clustering Relations into Abstract ER Schemas for Database Reverse Engineering", *Prof. of the 3th European conference on Software Maintenance and Reengineering*, Netherlands, March 03-05 1999, pp.169-176.
- [12] Halpin, Terry, *Conceptual Schema & Relational Database Design the 2nd edition*, WytLytPub, USA, 1999.
- [13] Han, Jiawei and Micheline Kamber, *Data Mining*:

Concepts and Techniques, Morgan Kaufmann Publishers,
USA, 2001.

Deriving Conceptual Schema from XML Databases

ORIGINALITY REPORT

94%

SIMILARITY INDEX

93%

INTERNET SOURCES

94%

PUBLICATIONS

9%

STUDENT PAPERS

PRIMARY SOURCES

1

fportfolio.petra.ac.id

Internet Source

93%

2

Oviliani Yenty Yuliana. "Deriving Conceptual Schema from XML Databases", 2009 First Asian Conference on Intelligent Information and Database Systems, 04/2009

Publication

1%

Exclude quotes On

Exclude bibliography On

Exclude matches < 1%