

# Algorithms for Minimal Dependency Set and Membership Based on XML Functional Dependency and Multi-valued Dependency

Zhongping Zhang<sup>1,2</sup>, Chunzhen Fang<sup>1</sup>

<sup>1</sup>The School of Information Science and Engineering, Yanshan University,  
Qinhuangdao, Hebei, 066004, China

<sup>2</sup>The Key Laboratory for Computer Virtual Technology and System Integration of Hebei Province,  
Qinhuangdao, Hebei, 066004, China

Email: zpzhang@ysu.edu.cn, fcz\_chunzhen@163.com

**Abstract**—As the XML functional dependency and multi-valued dependency impact on the normalization design of semi-structured data, definitions of XML functional dependency, XML multi-valued dependency, path dependency base and the minimal dependency set are given in this paper. Algorithms for minimal dependency set and membership with path expression based on the coexistence of XFD and XMVD are then proposed. Finally, the correctness and termination of these algorithms are proved, and their time complexities are analyzed as well.

**Index Terms**—XML functional dependency, XML multi-valued dependency, path dependency base, membership, minimal dependence set

## I. INTRODUCTION

XML (eXtensible Markup Language) inherits the powerful function of SGML(Standard Generalized Markup Language)[1] and makes up the deficiencies of HTML. It is a standard which is used for data expression and data exchange on the Internet, providing an effective means for data description and data exchange on the Internet applications. It is widely used in many fields [2-5]. XML schema is an important concept in the field of XML and it is the first step to build database applications. Currently, DTD develops well, and it is widely used in the practical applications of the XML document. However, because of some unusual data dependencies in XML database, it may result in data redundancies and abnormal operations due to design flaws for DTD [6].

In recent years, scholars have done a lot of exploration and research on the normalization of XML database: the normalization based on functional dependency [7-12] and multi-valued dependency [13-15]. These research literatures analyze the effect on XML data and free

redundancy from a single point of view such as XML functional dependency or XML multi-valued dependency. However, they do not consider the effect on XML data and XML database normalization on the condition of coexistence of XML functional dependency and multi-valued dependency. Therefore, according to the inference rules based on the coexistence of XML functional dependency and XML multi-valued dependency [16], we propose algorithms for minimal dependency set (DEP-MINIMIZE algorithm) and membership (DEP-MEMBERSHIP algorithm) with path expression based on coexistence of XML functional dependency and multi-valued dependency. It simplifies the dependency set and ensures the simplification on analysis and calculation.

## II. PRELIMINARY DEFINITIONS AND NOTATIONS

In this section, we present some preliminary definitions and notations that we need.

**Definition 1:** (XML Tree) An XML tree is defined as  $T=(V,lab,ele,att,val,root)$ , it is said to conform to a DTD[9]  $D=(E_1,E_2,A,P,R,r)$ , denoted by  $T=D$ [17], where

- (i)  $T$  is the XML tree's name;
- (ii)  $V$  is a finite set of nodes in  $T$ ;
- (iii)  $lab$  is a function from  $V$  to  $E_1 \cup E_2 \cup A$ , which assigns a identifier to each node in  $V$ . A node  $v$  in  $V$  is called a complex element node if  $lab(v) \in E_1$ ; a simple element node if  $lab(v) \in E_2$ , and an attribute node if  $lab(v) \in A$ ;
- (iv)  $ele$  is a function from  $V$  to a sequence of  $V$  nodes, so that for any  $v \in V$ , if  $lab(v) \in E_1$ ,  $ele(v)$  is a set of some children of  $v$ . The node in  $ele(v)$  is element node, and if  $ele(v) = \{v_1, \dots, v_m\}$ , then  $\{lab(v_1), \dots, lab(v_m)\} \in P(lab(v))$ ;
- (v)  $att$  is a function from  $V$  to  $A$ . If  $att(v, l) = v_l$ , then  $lab(v) \in E_1$  and  $lab(v_l) = l$ ; if  $att(v) = \{v_1, \dots, v_n\}$ , then  $\{lab(v_1), \dots, lab(v_n)\} \in R(lab(v))$ , where  $v \in V, l \in A$ ;

Manuscript received December, 1, 2013; revised March 24, 2014; accepted April 21, 2014.

(vi) val is a function that assigns a value to each node. If a node  $v$  is a leaf node or a simple element node of  $T$ ,  $val(v)$  is a string value which is either the content of a text element or the content of an attribute; otherwise  $val(v)$  is the node's identifier of  $v$ .

(vii) root is the unique root node labeled with complex element name  $r$ .

Example 1: Describe a college's relationship among three entities Course, Student and Teacher and store data information in Relational database. Relational schema  $R$  of database is designed as following:

Course(Cno, Cname) Student(Sno, Sname)  
Teacher(Tno, Tname)  
Courses(Cno,Sno, Tno,Cname, Sname, Tname)

If the relational database need to be stored as an XML document, this document's DTD  $D$  is expressed as:

```
<!ELEMENT Courses (Course*)>
<!ELEMENT Course (Cname,Student*)>
  <!ATTLIST Course
    Cno CDATA #REQUIRED>
  <!ELEMENT Cname (#PCDATA)>
  <!ELEMENT Student (Sname,Teacher)>
    <!ATTLIST Student
      Sno CDATA #REQUIRED>
    <!ELEMENT Sname (#PCDATA)>
    <!ELEMENT Teacher (Tname)>
      <!ATTLIST Teacher
        Tno CDATA #REQUIRED>
    <!ELEMENT Tname (#PCDATA)>
```

There are some data instances in TABLE I:

TABLE I.  
A PART OF STUDENTS' COURSE RECORDS OF ONELEGE COL

Cno	Cname	Sno	Sname	Tno	Tname
C01	XML	S001	Amily	T001	Mary
C01	XML	S002	Stephen	T001	Mary
C01	XML	S001	Amily	T002	Anne
C01	XML	S002	Stephen	T002	Anne

According to this course records, we can obtain an XML tree  $T$  which based on the DTD  $D$ . It is shown in Figure 1.

Definition2: (Path Instance on XML Tree) Let  $T$  be an XML tree that satisfied the given DTD  $D$ . A path instance [7] over an XML tree  $T$  is a sequence,  $v_1 = v_r$  and for every  $v_i, 2 \leq i \leq n, v_i \in V$  and  $v_i$  is a child of  $v_{i-1}$ . A path  $v_1.v_2. \dots.v_{n-1}.v_n$  is defined as one path instance based on the path  $p_1.p_2. \dots.p_{n-1}.p_n$  if for all  $v_i, 1 \leq i \leq n, lab(v_i)=p_i$ . All path instances based on a path  $p$  over a tree  $T$  form a set which denoted by  $Paths(p)$ . All Path sets on  $D$  are defined as the  $Paths(D)$ .

Example 2: As Figure 1 shows that,  $v_1.v_2.v_4$  is a path instance of path  $Courses.Course.Cno$ ,  $Paths(Courses.Course.Cno)=\{v_1.v_2.v_4, v_1.v_3.v_6\}$ .

Definition 3:(XFD) Let  $T$  be an XML tree that satisfied the given DTD  $D$ . An XFD is a statement of the form  $p_1, \dots, p_k \rightarrow q_1, \dots, q_m, k \geq 1, m \geq 1, P=\{p_1, \dots, p_k\}$  and  $Q=\{q_1, \dots, q_m\}$  are subsets in  $Paths(D)$ . There are arbitrary two distinct path instances  $V=\{v_1^1.v_2^1. \dots.v_{l-1}^1.v_l^1, \dots, v_1^m.v_2^m. \dots.v_{n-1}^m.v_n^m\}$

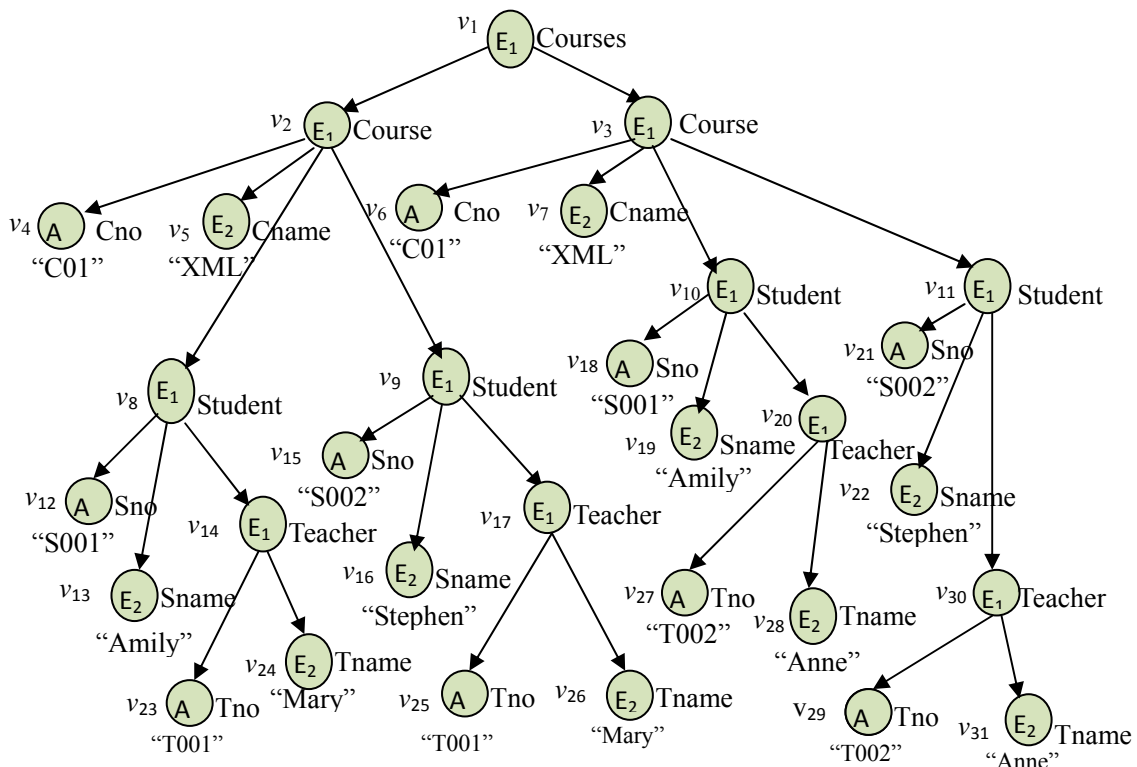


Figure 1. An XML tree formed from a part of students' course records

and  $W = \{ w_1^1 \cdot w_2^1 \cdot \dots \cdot w_{t-1}^1 \cdot w_t^1, \dots, w_1^m \cdot w_2^m \cdot \dots \cdot w_{n-1}^m \cdot w_n^m \}$  in  $\text{Paths}(Q)$ ,  $n \geq 1, t \geq 2$ ,  $T$  satisfies the XFD:  $p_1, \dots, p_k \rightarrow q_1, \dots, q_m$ , where

- (i)  $Q \subseteq P$ ; or
- (ii) For any two distinct path instances  $v_1^i \cdot v_2^i \cdot \dots \cdot v_{t-1}^i \cdot v_t^i$  and  $w_1^i \cdot w_2^i \cdot \dots \cdot w_{t-1}^i \cdot w_t^i$  in  $\text{Paths}(q_i)$ , if  $\text{Last}(p_j) [16] \in E_1$  and  $x_{ij} = y_{ij}$ , or  $\text{Last}(p_j) \notin E_1$  and  $\text{val}(\text{Nodes}(x_{ij}, p_j) [16]) \cap \text{val}(\text{Nodes}(y_{ij}, p_j)) \neq \emptyset$  then  $\text{val}(v_t^i) = \text{val}(w_t^i)$ ; where  $x_{ij} = \{ v | v \in \{ v_1^i, \dots, v_t^i \} \wedge v \in N(p_j \cap q_i) \}$ ,  $y_{ij} = \{ v | v \in \{ w_1^i, \dots, w_t^i \} \wedge v \in N(p_j \cap q_i) \}$ ,  $1 \leq j \leq k, 1 \leq i \leq m, t \geq 1$ .

We note that the path  $p_j \cap q_i$  is a prefix of  $q_i$ . There exists only one node in the set  $\{ v_1, \dots, v_t \}$  also in  $N(p_j \cap q_i)$ , therefore  $x_{ij}$  contains only one node.  $y_{ij}$  is similar to  $x_{ij}$  and it also contains only one node.

Definition 4: (XMVD) Let  $T$  be an XML tree that conforms to a DTD  $D$ . An XMVD is a statement of the form  $p_1, \dots, p_k \twoheadrightarrow q_1, \dots, q_m | r_1, \dots, r_s, 1 \leq k, 1 \leq m, 1 \leq s$ .  $P = \{ p_1, \dots, p_k \}$ ,  $Q = \{ q_1, \dots, q_m \}$  and  $R = \{ r_1, \dots, r_s \}$  are subsets in  $\text{Paths}(D)$ ,  $\{ \{ p_1, \dots, p_k \} \cup \{ q_1, \dots, q_m \} \cup \{ r_1, \dots, r_s \} \} \subseteq \text{Paths}(D)$ . The tree  $T$  satisfies the XMVD if there exists a  $q_i, 1 \leq i \leq m$ , and two distinct path instances  $v_1^i \cdot v_2^i \cdot \dots \cdot v_{t-1}^i \cdot v_t^i$  and  $w_1^i \cdot w_2^i \cdot \dots \cdot w_{t-1}^i \cdot w_t^i$  in  $\text{Paths}(q_i)$ ,  $1 \leq t$ , where

- (i)  $\text{val}(v_t^i) \neq \text{val}(w_t^i)$ ;
- (ii) There exists a  $r_j \in R, 1 \leq j \leq s$ , and two nodes  $z_1, z_2$ , where  $z_1 \in \text{Nodes}(x_{ij}, r_j)$  and  $z_2 \in \text{Nodes}(y_{ij}, r_j)$ , such that  $\text{val}(z_1) \neq \text{val}(z_2)$ ;
- (iii) For all  $p_h, 1 \leq h \leq k$ , there exists two nodes  $z_3$  and  $z_4$ , where  $z_3 \in \text{Nodes}(x_{ijh}, p_h)$  and  $z_4 \in \text{Nodes}(y_{ijh}, p_h)$ , such that  $\text{val}(z_3) = \text{val}(z_4)$ ;
- (iv) There exists a path instance  $v_1^i \cdot v_2^i \cdot \dots \cdot v_{t-1}^i \cdot v_t^i$  in  $\text{Paths}(q_i)$ , we have  $\text{val}(v_t^i) = \text{val}(w_t^i)$ , and there is a node  $z'_1$  in  $\text{Nodes}(x'_{ij}, r_j)$ , such that  $\text{val}(z'_1) = \text{val}(z_2)$ . There exists a node  $z'_3$  in  $\text{Nodes}(x'_{ijh}, p_h)$  such that  $\text{val}(z'_3) = \text{val}(z_3)$ ;
- (v) There is a path instance  $w_1^i \cdot w_2^i \cdot \dots \cdot w_{t-1}^i \cdot w_t^i$  in  $\text{Paths}(q_i)$  making  $\text{val}(w_t^i) = \text{val}(v_t^i)$ , and there exists a node  $z'_2$  in  $\text{Nodes}(y'_{ij}, r_j)$ , we have  $\text{val}(z'_2) = \text{val}(z_1)$  and there exists a node  $z'_4$  in  $\text{Nodes}(y'_{ijh}, p_h)$  such that  $\text{val}(z'_4) = \text{val}(z_4)$ .

Where,  $x_{ij} = \{ v | v \in \{ v_1^i, \dots, v_t^i \} \text{ and } v \in N(r_j \cap q_i) \}$ ,  $y_{ij} = \{ v | v \in \{ w_1^i, \dots, w_t^i \} \text{ and } v \in N(r_j \cap q_i) \}$ ,  $x'_{ij} = \{ v | v \in \{ v_1^i, \dots, v_t^i \} \text{ and } v \in N(p_h \cap r_j \cap q_i) \}$ ,  $y'_{ij} = \{ v | v \in \{ w_1^i, \dots, w_t^i \} \text{ and } v \in N(p_h \cap r_j \cap q_i) \}$ ;  $x'_{ijh} = \{ v | v \in \{ v_1^i, \dots, v_t^i \} \text{ and } v \in N(r_j \cap q_i) \}$ ,  $y'_{ijh} = \{ v | v \in \{ w_1^i, \dots, w_t^i \} \text{ and } v \in N(r_j \cap q_i) \}$ ,  $x'_{ijh} = \{ v | v \in \{ v_1^i, \dots, v_t^i \} \text{ and } v \in N(r_j \cap q_i) \}$ ,  $y'_{ijh} = \{ v | v \in \{ w_1^i, \dots, w_t^i \} \text{ and } v \in N(r_j \cap q_i) \}$ .

$= \{ v | v \in \{ v_1^i, \dots, v_t^i \} \text{ and } v \in N(p_h \cap r_j \cap q_i) \}$ ,  $y'_{ijh} = \{ v | v \in \{ w_1^i, \dots, w_t^i \} \text{ and } v \in N(p_h \cap r_j \cap q_i) \}$ .

We note that the path  $r_j \cap q_i$  is a prefix of  $q_i$ , there exists only one node in set  $\{ v_1^i, \dots, v_t^i \}$  also in  $N(r_j \cap q_i)$ , therefore  $x_{ij}$  contains only one node.  $y_{ij}, x'_{ijh}, y'_{ijh}, x'_{ij}, y'_{ij}, x'_{ijh}, y'_{ijh}$  are similar to  $x_{ij}$ . The XMVD is symmetrical, i.e. the XMVD:  $p_1, \dots, p_k \twoheadrightarrow q_1, \dots, q_m | r_1, \dots, r_s$  holds iff the XMVD:  $p_1, \dots, p_k \twoheadrightarrow r_1, \dots, r_s | q_1, \dots, q_m$ .

Definition 5: (Minimal Base) Let  $T$  be an XML tree that satisfied the given DTD  $D$ , Path set is  $P = \{ P_1, \dots, P_k \}$ . We assume that  $\text{Paths}(D) = P_1 \cup \dots \cup P_k$ . The minimal base of  $P$  is denoted by  $\text{MB}(P)$  and it is a partition of  $\text{Paths}(D)$ ,  $S_1, \dots, S_q$ , where

- (i) Each  $P_i$  is a set of  $S_j$  by union operation;
- (ii) There exists no partition that satisfies the condition (i) and the number of partition is less than  $q$ , where  $1 \leq i \leq k, 1 \leq j \leq q$ .

Definition 6: (Dependency Base) Let  $\Sigma$  be the set of XFD and XMVD over the complete instance document XML tree  $T$  which satisfied DTD  $D$ ,  $P \subseteq \text{Paths}(D)$ , the minimal base  $\text{MB}(P^+)$  of  $P^+$  [16] is a path dependency base relative to  $\Sigma$ , denoted by  $\text{DEP}(P)$ .

Definition 7: (Logical Implication) Let  $\text{Paths}(D)$  be the path set of DTD  $D$ , and  $\Sigma$  is the data dependency set of  $D$ . If each XML tree  $T$  of  $D$  satisfies  $\Sigma$  and  $P \twoheadrightarrow Q | R$  or  $P \rightarrow Q$ , we call  $P \twoheadrightarrow Q | R$  or  $P \rightarrow Q$  implicated logically by  $\Sigma$ , and denotes  $\Sigma \models P \twoheadrightarrow Q | R$  or  $\Sigma \models P \rightarrow Q$ .

### III. MEMBERSHIP ALGORITHM ON THE CONDITION OF COEXISTENCE OF XML FUNCTIONAL DEPENDENCY AND XML MULTI-VALUED DEPENDENCY

#### A.1. Path Dependency Base Algorithm

Dependency base is a partition of attribute set of relational data schema in the relational data theory. The path dependency base can gain the logical implication of multi-valued dependency directly; in other words, when the dependency base based on given multi-valued dependency is confirmed, we can get all multi-valued dependencies implicated logically by some attribute set.

Lemma 1: Let  $T$  be an XML tree that satisfied the given DTD  $D$ .  $p_h, r_j, q_i \in D, 1 \leq h \leq k, 1 \leq i \leq m, 1 \leq j \leq s$ . If XML tree  $T$  satisfies XFD:  $p_1, \dots, p_k \rightarrow q_1, \dots, q_m$ ,  $T$  will satisfy XMVD:  $p_1, \dots, p_k \twoheadrightarrow q_1, \dots, q_m | r_1, \dots, r_s$ , where  $R = \{ r_1, \dots, r_s \} \in D$ .

#### B.1. Algorithm Description

We note that if there is one dependency: XMVD  $P \twoheadrightarrow Q$  (XFD  $P \rightarrow Q$ ), and dependency set  $\Sigma$  implicates logically this dependency, so  $Q$  is the union set of some paths among  $\text{DEP}(P)$  according to the definition 6 and 7. The algorithm for path dependency base is given before solving the membership on the condition of coexistence of XML functional dependency and XML multi-valued dependency:

First, extend all the XFDs among dependency set  $\Sigma$  to XMVD, then we can get a transformation from  $\Sigma$  to  $\Sigma'$  which only includes XMVD; the initial value of BASIS

consists of all the single paths from path set  $P$  and  $\text{Paths}(D)-P$ . The initial value of change-flag is set as T, then execute “while” loop: initialize the value of change-flag as F, and then do the following operations for each MVD in  $\Sigma'$ : let  $P'$  be the non-empty set, all of its elements from the intersection of BASIS and  $P$ ,  $Q'=Q-P'$ . At the first time,  $W$  is empty. For the  $g_i$  in BASIS, add  $g_i$  into  $W$  if  $g_i \subset Q'$ . If  $Q' \neq \emptyset$  and  $Q' \neq W$ , set change-flag as T, then add  $Q'$  into BASIS by the definition 5. Executing the “while” loop repeatedly before  $Q' = \emptyset$  or  $Q'$  has been the union set of some elements from BASIS.

Algorithm 1: DEP-BASE(Path Dependency Base).

INPUT: mixed set  $\Sigma$  consisting of XFD and XMVD,  $\text{Paths}(D)$ ,  $P = \{p_1, \dots, p_k\}$ ;  
 OUTPUT: DEP( $P$ ).  
 DEP-BASE( $\Sigma, \text{Paths}(D), P$ )  
 Begin  
 (1) To transform  $\Sigma$  into  $\Sigma'$  including XMVD only.  
 (2) BASIS :=  $\{\{p_i\} \mid p_i \in P\} \cup \{\text{Paths}(D)-P\}$ ;  
 (3) change-flag := ‘T’;  
 (4) while change-flag do  
 (5) {change-flag := ‘F’;  
 (6) for every MVD  $P \rightarrow Q \in \Sigma'$  do  
 (7)  $\{P' := \cup \{R \mid R \in \text{BASIS and } R \cap P \neq \emptyset\}$ ;  
 (8)  $Q' := Q - P'$ ;  
 (9)  $W := \emptyset$ ;  
 (10) for  $g_i \in \text{BASIS}$  do  
 (11) {if  $g_i \subset Q'$  then  
 (12)  $\{W := W \cup g_i\}$ ;  
 (13) if  $Q' \neq \emptyset$  and  $Q' \neq W$  then do  
 (14) {change-flag := ‘T’;  
 (15) BASIS := MB{BASIS  $\cup$   $Q'$ };}  
 (16) return(BASIS)  
 End

### C.1. Analysis of Algorithm

#### C.1.1. Termination of the Algorithm

After the initialization of step (2), step (3) and the step (4)’s loop operation, BASIS becomes a partition of  $\text{Paths}(D)$ , because every subset after partitioning is non-empty, and the number of sets attained by partitioning the  $\text{Paths}(D)$  is at most the number of all the paths in  $\text{Paths}(D)$ . After every time (except the last time) for executing step (4), the size of BASIS will always increase. We note that the size of BASIS is at most the number of all the Paths in  $\text{Paths}(D)$ . Because the XMVD in  $\Sigma'$  is finite, so that the “for” loop in step (6) is terminable. In conclusion, algorithm 1 is terminable.

#### C.1.2. Correctness of the Algorithm

Correctness of the algorithm is to prove that BASIS is equal to DEP( $P$ ) when the algorithm terminates. The proof process includes two parts. The first part is to prove  $\text{BASIS} \subset \text{DEP}(P)$ , and the second part is to prove  $\text{DEP}(P) \subset \text{BASIS}$ .

First, to prove  $\text{BASIS} \subset \text{DEP}(P)$ . When the algorithm terminates,  $\text{BASIS} = \{\{P_1\}, \dots, \{P_m\}, \{p_1\}, \dots, \{p_k\}\}$ , where  $\{P_1\}, \dots, \{P_m\}$  is a partition of  $\text{Paths}(D)-P$ , and single paths of  $P$  form a set only including one path. According to the

reflexivity inference rules, we note that  $\Sigma$  implicates logically XMVD  $P \rightarrow p_i$ . The induction method is used to prove  $P \rightarrow P_j \in \Sigma^+$ , for each  $\{P_j\} \in \text{BASIS}$ , where  $1 \leq j \leq m$  as following.

According to the reflexivity inference rules, we have  $P \rightarrow Q \in \Sigma^+$ , so that  $P \rightarrow \forall \text{Paths}(D)-P-Q \in \Sigma^+$  holds. After executing the step(2), all elements in BASIS depend on  $P$  by multi-value, that is to say  $P \rightarrow P_j \in \Sigma^+$  hold, for each  $\{P_j\} \in \text{BASIS}$ , where  $1 \leq j \leq m$ .

After  $t(t \geq 0)$  times loop, we assume  $P \rightarrow P_j \in \Sigma^+$  hold, for each  $\{P_j\} \in \text{BASIS}$ , where  $1 \leq j \leq m$ . For the  $(t+1)^{\text{th}}$  loop: if  $P \rightarrow Q \in \Sigma$  is an XMVD during the  $(t+1)^{\text{th}}$  loop, BASIS value will change (if BASIS values don’t change, it will be the final value). Let  $P'$  be the non-empty set whose elements are from the intersection of BASIS and  $P$ .  $P' \rightarrow Q$  holds in accordance with the augmentation inference rules, because BASIS is a partition of  $\text{Paths}(D)$ ,  $P \subset P'$ .  $P \rightarrow P_j \in \Sigma^+$  holds, for each  $\{P_j\} \in \text{BASIS}$ , where  $1 \leq j \leq m$  before the  $(t+1)^{\text{th}}$  loop, then  $P \rightarrow P'$  holds under the union rules. At the same time, we note that  $P \rightarrow Q - P' \in \Sigma^+$  holds according to transitivity rules.  $Q' = Q - P'$  is set. If  $Q'$  is empty or  $Q'$  is a union set of some elements among BASIS,  $Q'$  is added into BASIS. Then solve the minimal base of BASIS. BASIS value doesn’t change according to definition 6. We need to make a corresponding modification for BASIS when  $Q'$  is not empty and  $Q'$  isn’t a union set of some elements among BASIS. Now, according to the difference rules,  $P \rightarrow P_j$  holds after modifying, for each  $\{P_j\} \in \text{BASIS}$ . In conclusion,  $P \rightarrow P_j \in \Sigma^+$  holds, for every  $\{P_j\} \in \text{BASIS}$ ,  $1 \leq j \leq m$  during the  $(t+1)^{\text{th}}$  loop, so that  $\text{BASIS} \subset \text{DEP}(P)$  holds in accordance with the definition 6.

Second, to prove  $\text{DEP}(P) \subset \text{BASIS}$ . Constructing a XML tree  $T$  that conforms to DTD  $D$ , where

(i) Each XMVD on  $T$  from  $\Sigma$  is legitimate;

(ii) The necessary and sufficient condition of one XMVD  $P \rightarrow Q$  on  $T$  holds is that  $Q'$  is a union set of some elements among BASIS.

Constructing an XML document tree  $T$ : there are  $2^m$  tree tuples, and every tuple in  $\text{Tuple}_T(D)$  [16] has a group of corresponding sequence  $\{a_1, \dots, a_m\}$ , where  $a_i \in [0, 1]$ .  $\{P_1\}, \dots, \{P_m\}$  is a partition of  $\text{Paths}(D)-P$ . The corresponding values of all paths in  $P^+$  for every tuple in  $\text{Tuple}_T(D)$  are 1, and the value of every path in  $P_i$  is  $a_i$ .

Properties of the XML document tree  $T$ :

Property 1: Each XMVD is valid when on the right of  $T$  is  $P_i$ .

Property 2: Some XMVD holds on  $T$  when on the right of  $P_i$  is a non-empty subset iff the left of XMVD has intersection with  $P_i$ .

Each multi-valued dependency on the  $T$  will be proved correct as following:

$P \rightarrow Q \in \Sigma$  holds, and  $P'$  is the set whose elements attained from the intersection of BASIS and  $P$ . We note that  $Q - P'$  is an empty set or the union set of some elements among BASIS in accordance with termination of the algorithm, so that  $P \rightarrow Q - P'$  holds on  $T$ . According to Property 2,  $P \rightarrow Q \cap P'$  holds on  $T$ , then  $W \rightarrow R$  holds on  $T$  based on union rules among

deriving.

We will prove  $P \twoheadrightarrow Q$  hold on  $T$  iff  $Q$  is the union set of some elements in BASIS.

We note from definition 5: if  $Q$  is the union set of some elements in BASIS,  $P \twoheadrightarrow Q$  holds on  $T$ , that is to say the sufficiency condition holds.

If  $P \twoheadrightarrow Q$  holds on  $T$ ,  $P \twoheadrightarrow Q \cap P_i$  holds on  $T$ , for each  $i$ ,  $1 \leq i \leq m$ . Because the intersection of  $P$  and  $P_i$  is empty,  $Q \cap P_i$  is empty or is  $P_i$  based on Property 1 and Property 2. Therefore,  $Q$  is the union set of some elements among BASIS.

In conclusion,  $\text{BASIS} = \text{DEP}(P)$ , algorithm 1 is correct.

### C.1.3. Time Complexity of the Algorithm

“ $m$ ” expresses the total number of paths from  $\text{Paths}(D)$ , and “ $n$ ” expresses the number of dependencies from  $\Sigma$ .

Because XFD is a special case of XMVD, so we note that every XFD can be extended to the corresponding XMVD.

Partition the path set into  $g$  subsets, using matrix  $g \times m$  for expressing. Each row with  $m$  digits expresses a set. Each column of the matrix has only one 1, other positions are 0. The  $(ij)^{\text{th}}$  position of the matrix is 1 iff path  $p_j$  belongs to the  $i^{\text{th}}$  set. Obviously,  $p_j$  and  $p_i$  are in the same path set iff the  $i^{\text{th}}$  column and the  $j^{\text{th}}$  column of the matrix are the same; in other words, when the set is on the  $h^{\text{th}}$  row, the values of  $h_i$  and  $h_j$  both are 1, then other positions of the  $i^{\text{th}}$  column and the  $j^{\text{th}}$  column of the matrix are 0. Therefore, to find the minimal path dependency base, we partition the same columns into one set whose row values are 1 in the same columns. Arbitrary two of all the columns are compared for  $g$  times. We note that the time cost of which solving the minimal path dependency base is  $O(g \times m)$ .

This algorithm initializes the path dependency base at first. Every single paths from path  $P$  set form a set respectively and all sets constitute one set, then path set  $\text{Paths}(D) - P$  forms one set. Since constituting two rows of the matrix after initialization, we get the time cost of the operation is  $O(m)$ .

The loop in step (4) executes at most  $m$  times. In every loop, for the given XMVD  $P \twoheadrightarrow Q$ , our goal is to find the union set of all the elements attained from the intersection of BASIS and  $Q$ . Since the size of BASIS is at most  $m$ , the time cost of the operation is  $O(m^2)$ . The time complexity of the “for” loop is  $O(m^2)$  of step (10). In step (4), it is possible to examine whether  $n$  dependencies change the BASIS value or not. Therefore, the loop continuous until the BASIS value change, and it need to execute  $O(n \times m^2)$  times. The time cost of solving minimal path dependency base is  $O(m^2)$  and the time cost of executing step(4) once is  $O(n \times m^2)$ . After executing step (4) completely, the total time cost is  $O(n \times m^3)$ .

In conclusion, the total time complexity of this algorithm is  $O(n \times m^3)$ .

### A.2. Membership Algorithm

The membership is to solve whether some dependency implicated logically by dependency set or not. On the basis of DEP-BASE algorithm, we provide the

membership algorithm in the following.

### B.2. Algorithm Description

First, we get the path dependency base, and the initial value of  $Q'$  is set as empty. If the dependency is a multi-valued dependency, the following operations is executed: if the element in  $\text{DEP}(P)$  belongs to the right path set of this multi-valued dependency, then add this element into  $Q'$ . Whether  $Q' = Q$  hold or not examined. If holds, this dependency is implicated logically by  $\Sigma$ . If the dependency is functional dependency, and need to examine whether the right path set  $Q \in P^+$  of this dependency hold or not; if holds, this dependency is implicated logically by  $\Sigma$ .

Algorithm 2: DEP-MEMBERSHIP(XFD and XMVD Membership).

INPUT: dependency set  $\Sigma$ , path set  $\text{Paths}(D)$ , a dependency  $g: P \twoheadrightarrow Q (P \rightarrow Q)$ .

OUTPUT: if this dependency is implicated logically by  $\Sigma$ , the output is True; otherwise, the output is False.

DEP-MEMBERSHIP( $\Sigma, g$ )

Begin

(1)  $\text{DEP}(P) := \text{DEP-BASE}(\Sigma, \text{Paths}(D), P)$ ;

(2)  $Q' := \emptyset$ ;

(3) if  $P \twoheadrightarrow Q$

(4) for every  $P_i$  in  $\text{DEP}(P)$  do

(5) if  $P_i \subseteq Q$  then

(6)  $Q' := Q' \cup P_i$ ;

(7) if  $(P \twoheadrightarrow Q \text{ and } Q' = Q)$  or  $(P \rightarrow Q \text{ and } Q \in P^+)$  then

(8) return(True);

(9) else

(10) return(False);

End

### C.2. Analysis of Algorithm

#### C.2.1. Termination Of the Algorithm

The step (1) of algorithm 2 calls the algorithm 1 to get the  $\text{DEP}(P)$ , so the step(1) is terminable by algorithm 1. Because the number of elements from  $\text{DEP}(P)$  is finite and at most the number of paths in  $\text{Paths}(D)$ , the “for” loop is terminable in step (4) of algorithm 2. Therefore, the algorithm 2 is terminable.

#### C.2.2. Correctness of the Algorithm

If the dependency needed to be examined is XMVD, it can be transformed from judging whether  $\Sigma$  implicates logically the XMVD  $P \twoheadrightarrow Q$  to judging whether the  $Q$  is the union set of some elements from  $\text{DEP}(P)$  by definition 6. The step (1) of algorithm 2 calls the algorithm 1 to attain the  $\text{DEP}(P)$ , and we note that the step (1) is correct. The “for” loop in step(4) is to examine one by one whether every set belongs to  $Q$ , then  $Q'$  is used for expressing the sets whose elements in  $\text{DEP}(P)$  also in  $Q$ . Finally, compare  $Q'$  with  $Q$ . If they are same,  $Q$  is the union set of some elements from  $\text{DEP}(P)$ , now return True, otherwise, return False. If the dependency needed to be examined is XFD, we note that algorithm 2 is correct by the definition of closure.

### C.2.3. Time Complexity of the Algorithm

The number of paths from Paths(D) is expressed as  $m$ , and  $n$  expresses the number of dependencies from  $\Sigma$ . The time complexity of the step (1) in this algorithm is  $O(n \times m^3)$ . The “for” loop in step (4) mainly make a comparing among sets, and the number of elements from DEP(P) is at most  $m$ , so its time complexity is  $O(m^2)$ . Therefore, the total time cost of algorithm 2 is  $O(n \times m^3)$ .

## IV. MINIMAL DEPENDENCY SET ALGORITHM

Definition 8: Let  $\Sigma$  be the set of XFD and XMVD, if  $\Sigma_{\min}$  is a minimal dependency set of  $\Sigma$ , the  $\Sigma_{\min}$  should meet the following conditions:

(i) There is no redundancy path in  $\Sigma_{\min}$ , that is to say there is no  $P \rightarrow \rightarrow Q$  (or  $P \rightarrow Q$ ) that makes

$$\Sigma_{\min} = \Sigma_{\min} - \{P \rightarrow \rightarrow Q\} \text{ (or } \Sigma_{\min} = \Sigma_{\min} - \{P \rightarrow Q\} \text{)}.$$

(ii) There is no redundancy on the left of every dependency, namely for every XMVD  $P \rightarrow \rightarrow Q$  (or XFD  $P \rightarrow Q$ )  $\in \Sigma_{\min}$ , there not exist XMVD  $P' \rightarrow \rightarrow Q$  (or XFD  $P' \rightarrow Q$ )  $\in \Sigma^+$  that makes  $P' \subset P$  hold.

(iii) There is no redundancy on the right of every dependency, and the right of every dependency is single path, namely for every XMVD  $P \rightarrow \rightarrow Q$  (or XFD  $P \rightarrow Q$ )  $\in \Sigma_{\min}$ , there not exist XMVD  $P \rightarrow \rightarrow Q'$  (or XFD  $P \rightarrow Q'$ )  $\in \Sigma^+$  that makes  $\emptyset \subset Q' \subset Q$  hold.

### A. Algorithm Description

Data dependency plays an important role in the normalization design of database [18]. Designing an XML database that can avoid data redundancy and abnormal operation is an important research subject in the XML field [19].

Some data dependency can be implicated logically by other data dependencies. Removing the redundancy dependency and redundancy path is to simplify the given dependency set and make sure it is simple during the analysis and computation, on the condition that the data dependency set closure doesn't change. This problem is about minimal dependency set.

First, we set change-flag as T, then execute “while” loop, let change-flag be F, and initialize the value of  $\Sigma$  as  $\Sigma'$ . To perform the following operations for every  $d$  in  $\Sigma$ : judging whether the dependency is redundancy using DEP-MEMBERSHIP (algorithm 2) at first. If it is redundant, then remove it from  $\Sigma$ . If it is not redundant then to examine whether the left and right of this dependency exist the redundancy paths respectively, and if there is a redundancy path, it will be removed. After executing the operations, we need to examine whether  $\Sigma \neq \Sigma'$  hold or not, if it holds,  $\Sigma$  is not the minimal dependency set. Then set the change-flag as T and re-execute the “while” loop. If  $\Sigma = \Sigma'$ ,  $\Sigma$  is a minimal dependency set.

Algorithm3: DEP-MINIMIZE(Minimize Dependency).

INPUT: A set  $\Sigma$  including XFD and XMVD.

OUTPUT: the minimal dependency set  $\Sigma_{\min}$  of  $\Sigma$ .

DEP-MINIMIZE( $\Sigma$ )

Begin

```

(1)change-flag:='T';
(2)while change-flag do
(3)change-flag:='F';
(4) $\Sigma' := \Sigma$ ;
(5)for (every dependency  $d$  in  $\Sigma$ ) do
(6) if DEP_MEMBERSHIP( $\Sigma-d, d$ ) then
(7)    $\Sigma := \Sigma-d$ ;
(8) if (exist redundancy path on the left of  $d$ )
(9)   for every path  $p \in P$  do
(10)    if  $d$  is a XFD, then  $P \rightarrow Q$  do
(11)     if DEP_MEMBERSHIP( $\Sigma, \{P-p\} \rightarrow Q$ )
        then
(12)       $\Sigma := (\Sigma - \{P \rightarrow Q\}) \cup \{P-p \rightarrow Q\}$ ;
(13) if  $d$  is a XMVD, then  $P \rightarrow \rightarrow Q$  do
(14)   if DEP_MEMBERSHIP( $\Sigma, \{P-p\} \rightarrow \rightarrow Q$ )
        then
(15)     $\Sigma := (\Sigma - \{P \rightarrow \rightarrow Q\}) \cup \{P - \{p\} \rightarrow \rightarrow Q\}$ ;
(16) if (exist redundancy attribute on the right of  $d$ )
(17)   for every path  $q \in Q$  do
(18)    if  $d$  is a XFD, then  $P \rightarrow Q$  do
(19)      $\Sigma := (\Sigma - \{P \rightarrow Q\}) \cup \{P \rightarrow q_1, \dots, P \rightarrow q_m\}$ ;
(20) if  $d$  is a XMVD, then  $P \rightarrow \rightarrow Q$  do
(21)   if DEP_MEMBERSHIP( $\Sigma, P \rightarrow \rightarrow \{Q\} - \{q\}$ )
        then
(22)  $\Sigma := (\Sigma - \{P \rightarrow \rightarrow Q\}) \cup \{P \rightarrow \rightarrow Q', P \rightarrow \rightarrow \{Q\} - \{q\}\}$ 
(23)if  $\Sigma \neq \Sigma'$ 
(24) change-flag:='T';
(25)return( $\Sigma$ ).
End

```

### B. Analysis of Algorithm

#### B.1. Termination of the Algorithm

The number of elements from preliminary  $\Sigma$  and from  $\Sigma^+$  is finite. In addition, the size of  $\Sigma$  is at most the size of  $\Sigma^+$  after executing the step (5). The size of  $\Sigma^+$  is finite which results in the “for” loop in step (5) is terminable. The left and right paths of every dependency are finite. We note that the “for” loop is terminable in step (9) and step (17) of algorithm 3. In conclusion, the algorithm 3 is terminable.

#### B.2. Correctness of the Algorithm

Examine whether every dependency conforms to the definition of minimal dependency set or not in “for” loop of step (5) after executing the step (1)~step (4) of this algorithm. step (6) and step(7) of algorithm calls DEP\_MEMBERSHIP algorithm to remove the redundancy dependency; step(8)~step(15) calls DEP\_MEMBERSHIP algorithm to remove the left redundancy path; step(18) and step (19) transforms the right path of XFD to single path; step (8)~step(15) calls DEP\_MEMBERSHIP algorithm to remove the right redundancy path for XMVD. We note that algorithm 3 is correct in that the dependency meets the definition of minimal dependency set.

#### B.3. Time Complexity of the Algorithm

“ $m$ ” expresses the total number of paths from Paths(D), and “ $n$ ” expresses the number of dependencies from  $\Sigma$ .

The time cost of step (6) is  $O(n \times m^3)$  based on calling the DEP\_MEMBERSHIP algorithm; the time cost of step (8)~step(15) is equal to the time cost of the “for” loop, and it is  $O(n \times m^4)$ . The time complexity of step(16)~step(22) also depends on the “for” loop whose time complexity is  $O(n \times m^4)$ . In addition, the number of dependencies  $\Sigma$  is  $n$ . Therefore, the time complexity of executing the step(5) once is  $O(n^2 \times m^4)$ . In conclusion, the time complexity of algorithm 3 is  $O(n^2 \times m^4)$ .

## V. CONCLUSION

The design of database schema is the first step of the database application. If the database schema is well designed, abnormal data dependencies, data redundancies and abnormal operations can be avoided. This paper studies the membership problem from the view of the condition of coexistence of XML functional dependency and XML multi-valued dependency. At the same time, we propose DEP-BASE algorithm, DEP-MEMBERSHIP algorithm and DEP-MINIMIZE algorithm on the membership problem, not only proving the termination and correctness of these algorithms, but also analyzing their time complexities. They laid a foundation for designing a normal form which level of normalization is higher than others.

## ACKNOWLEDGMENT

This work was financially supported by Hebei Provincial Natural Science Foundation of China (F2012203087), National Natural Science Foundation of China (61272124) and National Natural Science Foundation of China (61073060).

## REFERENCES

- [1] Erik Naggam, “Standard Generalized Markup Language”[EB/OL].(1996-03-01)[2013-11-26] <http://www.w3.org/MarkUp/SGML>.
- [2] Arash Termehchy, Marianne Winslett, “Using Structural Information in XML Keyword Search Effectively”. ACM Transactions on Database Systems, 2011, vol.36, no.1, pp.4.
- [3] Haiping Xu, Abhinay Reddyreddy, Daniel F. Fitch, “Defending Against XML-Based Attacks Using State-Based XML Firewall”. Journal of Computers, 2011, vol. 6, no. 11, pp.2395-2407.
- [4] Ren Li, Jianhua Luo, Dan Yang, Haibo Hu, Ling Chen, “A Scalable XSLT Processing Framework based on MapReduce”. Journal of Computers, 2013, vol. 8, no. 9, pp.2175-2181.
- [5] Shihan Yang, Jinzhao Wu, Anping He, Yunbo Rao, “Derivation of OWL Ontology from XML Documents by Formal Semantic Modeling”. Journal of Computers, 2013, vol. 8, no. 2, pp. 372-379.
- [6] Haiyan Huang, Ronghua Shi, Gaoshi Li, “The normalization of the algorithm on XML multi-valued dependency”. Journal of hunan institute of science and technology, 2008, vol.29, no.12, pp. 126-129.
- [7] Millist W.Vincent, JIXUE LIU, CHENGFEI LIU, “Strong Functional Dependencies and Their Application to Normal Forms in XML”. ACM Transactions on Database Systems, 2004, vol.29, no.23, pp. 445-462.
- [8] E. F. Codd, “Recent investigations in relational data base systems”[C]//Proceedings of IFIP Congress 74, Stockholm, Sweden, 1974, pp.1017-1021.
- [9] Kamsuriah Ahmad, Ali Mamat, Hamidah Ibrahim, Shahrul Azman Mohd Noah, “Defining Functional Dependency for XML”. Journal of Information Systems, Research & Practices, 2008, vol.1, no.1, pp. 26-34.
- [10] Xiangguo Zhao, Junchang Xin, Ende Zhang, “XML Functional Dependency and Schema Normalization”[C]//Proceedings of the 9th International Conference on Hybrid Intelligent Systems, Shenyang, China, 2009, pp.307-312.
- [11] Marcelo Arenas, “Normalization Theory for XML”. SIGMOD Record, 2006, vol.35, no.4, pp. 57-64.
- [12] Tadeusz Pankowski, Tomasz Pilka, “Transformation of XML Data into XML Normal Form”. Informatica (Slovenia). 2009, vol.33, no.4, pp.417-430.
- [13] CATRIEL BEERI, “ON the Membership Problem for Functional and Multivalued Dependencies in Relational Database”. ACM Transactions on Database Systems. 1980, vol.5, no.3, pp.241-259.
- [14] Zhongping Zhang, “The logical implication algorithm research of XML multi-valued dependency”. Computer Science, 2006, vol.33, no. 11(Supplement), pp.353-354.
- [15] Wei Qiu, Lichen Zhang, “Study of Normalization Existing MVD in XML DTD”. Computer Science, 2007, vol.34, no.2, pp. 149-185.
- [16] Zhixiao Liu, “XML normalization research based on functional dependency and multi-valued dependency”, yanshan university, qinhuangdao, MA, 2012.
- [17] Millist W.Vincent, Jixue Liu, Chengfei Liu, Mukesh Mohania, “Multivalued Dependencies and a 4NF for XML”. CAISE 2003, pp. 14-29.
- [18] Jixue Liu, Jiuyong Li, Chengfei Liu, Yongfeng Chen., “Discover Dependencies from Data—A Review”. IEEE Transactions on Knowledge and Data Engineering, 2012, vol. 24, no.2, pp.251-264.
- [19] M. W. Vincent, J. Liu, M. Mohania, “The implication problem for ‘closest node’ functional dependencies in complete XML documents”. Journal of Computer and System Sciences, 2012, vol.78, no. 4, pp. 1045-1098.



**Zhongping Zhang**, Male, Born in 1972, professor, Ph.D., post-doctoral, CCF Senior Member (E20-0006458S). His main research interests are the grid computing, data mining and semi-structured data etc. He has undertaken 1 project of provincial level and has participated 2 projects funded by national natural science foundation of China. He rewarded the provincial scientific and technological progress second-class Award. On the domestic and international academic conferences and journals, He published more than 80 papers, 15 of them were cited by EI.

**Chunzhen Fang**, Female, Born in 1987, Postgraduate student, the main research interest is the outlier detection in data mining.