course: **Database Systems** (A7B36DBS)

lecture 9:

Relational design – algorithms

Doc. RNDr. Irena Holubova, Ph.D.

Acknowledgement: The slides were kindly lent by Doc. RNDr. Tomas Skopal, Ph.D., Department of Software Engineering, Charles University in Prague

Today's lecture outline

schema analysis

- basic algorithms (attribute closure, FD membership and redundancy)
- determining the keys
- testing normal forms
- normalization of universal schema
 - decomposition (to BCNF)
 - synthesis (to 3NF)

Attribute closure

functional dependency

- closure X⁺ of attribute set X according to FD set F
 - principle: we iteratively derive all attributes "F-determined" by attributes in X
 - complexity O(m*n), where n is the number of attributes and m is number of FDs

```
algorithm AttributeClosure(set of dependencies F, set of attributes X) :
  returns set X<sup>+</sup>
  ClosureX := X; DONE := false; m = |F|;
  while not DONE do
    DONE := true; Left-hand side of FD
    for i := 1 to m do
        if (LS[i] ⊆ ClosureX and RS[i] ⊈ ClosureX) then
        ClosureX := ClosureX ∪ RS[i];
        DONE := false;
    endif
    endfor
  endwhile
  return ClosureX;
```

The trivial FD is used (algorithm initialization) and then transitivity (test of left-hand side in the closure). The composition and decomposition usage is hidden in the inclusion test.

Example – attribute closure

$$F = \{a \rightarrow b, bc \rightarrow d, bd \rightarrow a\}$$

{b,c}⁺ = ?

1. ClosureX := {b,c} (initialization)2. ClosureX := ClosureX \cup {d} = {b,c,d} (bc \rightarrow d)3. ClosureX := ClosureX \cup {a} = {a,b,c,d} (bd \rightarrow a)

 $\{b,c\}^+ = \{a,b,c,d\}$

Membership test

- we often need to check <u>if a FD X → Y belongs to F</u>⁺, i.e., to solve the problem {X → Y} ∈ F⁺
- materializing F⁺ is not practical, we can employ the attribute closure

algorithm IsDependencyInClosure (set of FDs F,

FD X \rightarrow Y)

```
return Y \subseteq AttributeClosure(F, X);
```

Redundancy testing

The membership test can be easily used when testing redundancy of

- FD X \rightarrow Y in F
- attribute a in X (according to F and X \rightarrow Y)

```
algorithm IsDependencyRedundant (set of FDs F, FD X \rightarrow Y \in F)
return IsDependencyInClosure(F - {X \rightarrow Y}, X \rightarrow Y);
```

algorithm **IsAttributeRedundant**(set of FDs F, FD X \rightarrow Y \in F, attr. a \in X) **return** *IsDependencyInClosure*(F, X - {a} \rightarrow Y);

In the following slides we find useful the algorithm for <u>reduction</u> of the left-hand side of a FD:

Minimal cover

• for all FDs we test redundancies and remove them

```
algorithm GetMinimumCover(set of dependencies F)
  : returns minimal cover G
   decompose each FD in F into elementary FDs
   for each X \rightarrow Y in F do
       F := (F -
                                       removing redundant attributes
                \{X \rightarrow Y\} \cup
                {ReduceAttributes(F, X \rightarrow Y) \rightarrow Y};
   endfor
                                                  removing redundant FDs
   for each X \rightarrow Y in F do
       if IsDependencyRedundant(F, X \rightarrow Y)
       then F := F - \{X \rightarrow Y\};
   endfor
return F;
```

Determining (first) key

- the algorithm for attribute redundancy testing could be used directly for determining a key
- redundant attributes are iteratively removed from left-hand side of trivial FD A \rightarrow A

```
algorithm GetFirstKey(set of deps. F, set of attributes A)
```

```
: returns a key K;
return ReduceAttributes(F, A \rightarrow A);
```

<u>Note:</u> Because multiple keys can exists, the algorithm finds only one of them.

Which one? It depends on the traversing of the attribute set within the algorithm ReduceAttributes.

Determining all keys, the principle

Let us have a schema S(A, F). Simplify *F* to minimal cover.



- 1. Find any key K (see the previous slide).
- 2. Take a FD $X \rightarrow y$ in F such that $y \in K$ or terminate if not exists (there is no other key).
- 3. Because $X \rightarrow y$ and $K \rightarrow A$, it transitively holds also $X\{K y\} \rightarrow A$, i.e., $X\{K y\}$ is super-key.
- 4. Reduce FD $X{K y} \rightarrow A$ so we obtain key K' on the left-hand side.

This key is surely different from **K** (we removed **y**).

5. If **K'** is not among the determined keys so far, we add it, declare **K**=**K'** and continue from step 2. Otherwise we finish.

Determining all keys, the algorithm

Formally: Lucchesi-Osborn algorithm

- having an already determined key, we search for equivalent sets of attributes, i.e., other keys
- NP-complete problem (theoretically exponential number of keys/FDs)

```
algorithm GetAllKeys (set of FDs F, set of attr. A)
   : returns set of all keys Keys;
   let all dependencies in F be non-trivial
   K := GetFirstKev(F, A);
   Keys := \{K\};
   for each K in Keys do
      for each X \rightarrow Y in F do
          if (Y \cap K \neq \emptyset and \neg \exists K' \in Keys : K' \subseteq (K \cup X) - Y) then
              N := ReduceAttributes (F, ((K \cup X) - Y) \rightarrow A);
             Keys := Keys \cup {N};
          endif
      endfor
   endfor
return Keys;
```

Example – determining all keys

Contracts(A, F) $A = \{c = ContractId, s = SupplierId, j = ProjectId, d = DeptId, p = PartId, q = Quantity, v = Value\}$ $F = \{c \rightarrow all, sd \rightarrow p, p \rightarrow d, jp \rightarrow c, j \rightarrow s\}$

- 1. Determine the first key Keys = {C}
- 2. <u>Iteration 1:</u> take $jp \rightarrow c$ that has a part of the last key on the right-hand side (in this case the whole key c) and jp is not a super-set of already determined key
- 3. $jp \rightarrow all$ is reduced (no redundant attribute), i.e., Keys = {c, jp}
- 4. <u>Iteration 2:</u> take $sd \rightarrow p$ that has a part of the last key on the right-hand side (jp),

{jsd} is not a super-set of c nor jp, i.e., it is a key candidate

- 5. in jsd $\rightarrow all$ we get redundant attribute s, i.e., Keys = {c, jp, jd}
- 6. <u>Iteration 3:</u> take $p \rightarrow d$, however, jp was already found so we do not add it
- 7. Finish as the iteration 3 resulted in no key addition.

Testing normal forms

- NP-complete problem
 - we must know all keys then it is sufficient to test a FD in F, so we do not need to materialize F⁺
 - or, just one key needed, but also needing extension of F to F⁺
- fortunately, in practice determination of keys is fast
 - thanks to limited size of F and "separability" of FDs

Design of database schemas

Two ways of modelling a relational database:

- 1. we get <u>a set of relational schemas</u> (as either direct relational design or conversion from conceptual model)
 - normalization performed separately on each table
 - the database could get unnecessarily highly "granularized" (too many tables)
- 2. considering the whole database as a bag of (global) attributes results in <u>a</u> single universal database schema i.e., one big table + a single set of FDs
 - normalization performed on the universal schema
 - less tables (better "granulating")
 - "classes/entities" are generated (recognized) as the consequence of FD set
- both approaches could be combined i.e.,
 - create a conceptual database model
 - convert it to relational schemas
 - merge and/or normalize some of the schemas

Relational design – algorithms (A7B36DBS, Lect. 9)

Relational schema normalization

- just one way <u>decomposition</u> to multiple schemas
 - or merging some "abnormal" schemas and then decomposition
- different criteria
 - data integrity preservation
 - lossless join
 - dependency preserving
 - requirement on normal form (3NF or BCNF)
- manually or algorithmically

Why to preserve integrity?

If the decomposition is not limited, we can decompose the table into several single-column ones that surely are all in BCNF.

Company	HQ	Altitude		Company		HQ	Altitude
Sun	Santa Clara	25 m		Sun		Santa Clara	25 m
Oracle	Redwood	20 m		Oracle		Redwood	20 m
Microsoft	Redmond	10 m		Microsoft		Redmond	10 m
IBM	New York	15 m		IBM		New York	15 m
			-	Company	-	HQ	 Altitude

 $\frac{\text{Company}}{\text{HQ} \rightarrow \text{Altitude}}$

Clearly, there is something wrong with such a decomposition...

...it is **lossy** and it does not **preserve dependencies**

Lossless join

- a property of decomposition that ensures correct joining (reconstruction) of the universal relation from the decomposed ones
- <u>Definition 1</u>: Let $R({X \cup Y \cup Z}, F)$ be universal schema, where $Y \rightarrow Z \in F$. Then decomposition $R_1({Y \cup Z}, F_1)$, $R_2({Y \cup X}, F_2)$ is lossless.



• Alternative <u>Definition 2</u>: Decomposition of R(A, F) into $R_1(A_1, F_1)$, $R_2(A_2, F_2)$ is lossless, if $A_1 \cap A_2 \rightarrow A_1$ or $A_2 \cap A_1 \rightarrow A_2$



Example – lossy decomposition

Company	Uses DBMS	Data managed	
Sun	Oracle	50 TB	
Sun	DB2	10 GB	
Microsoft	MSSQL	30 TB	
Microsoft	Oracle	30 TB	

Company	Uses DBMS	Company	Data managed
Sun	Oracle	Sun	50 TB
Sun	DB2	Sun	10 GB
Microsoft	MSSQL	Microsoft	30 TB
Microsoft	Oracle	Company,	Data managed

Company, Uses DBMS

Company	Uses DBMS	Data managed
Sun	Oracle	50 TB
Sun	Oracle	10 GB
Sun	DB2	10 GB
Sun	DB2	50 TB
Microsoft	MSSQL	30 TB
Microsoft	Oracle	30 TB

Company, Uses DBMS



"reconstruction" (natural join)

Company, Uses DBMS, Data managed

Example – lossless decomposition

Company	HQ	Altitude
Sun	Santa Clara	25 m
Oracle	Redwood	20 m
Microsoft	Redmond	10 m
IBM	New York	15 m

	Company	HQ
	Sun	Santa Clara
,	Oracle	Redwood
	Microsoft	Redmond
	IBM	New York

HQ	Altitude
Santa Clara	25 m
Redwood	20 m
Redmond	10 m
New York	15 m

 $\frac{\text{Company}}{\text{HQ} \rightarrow \text{Altitude}}$

<u>Company</u>

<u>H0</u>

"reconstruction" (natural join)

Dependency preserving

- a decomposition property that ensures **no FD will be lost**
- Definition:

Let $R_1(A_1, F_1)$, $R_2(A_2, F_2)$ be decomposition of R(A, F). Then such decomposition preserves dependencies if $F^+ = (\bigcup_{i=1..n} F_i)^+$.

- Dependency preserving could be violated in two ways
 - during decomposition of F we do not derive all valid FDs we lose FD that should be preserved in a particular schema
 - even if we derive all valid FDs (i.e., we perform projection of F⁺), we may lose a FD that is valid across the schemas

Example – dependency preserving

dependencies not preserved, we lost $HQ \rightarrow Altitude$

Company	HQ	Altitude
Sun	Santa Clara	25 m
Oracle	Redwood	20 m
Microsoft	Redmond	10 m
IBM	New York	15 m

<u>Company</u>, HQ \rightarrow Altitude



Company	Altitude
Sun	25 m
Oracle	20 m
Microsoft	10 m
IBM	15 m

Company	HQ
Sun	Santa Clara
Oracle	Redwood
Microsoft	Redmond
IBM	New York

<u>Company</u>

<u>HQ</u>

Company	HQ
Sun	Santa Clara
Oracle	Redwood
Microsoft	Redmond
IBM	New York

HQ	Altitude
Santa Clara	25 m
Redwood	20 m
Redmond	10 m
New York	15 m

Company

<u>HQ</u>

The "Decomposition" algorithm

- algorithm for decomposition into BCNF, preserving lossless join
- may not preserve dependencies
 - not an algorithm property sometimes we simply cannot decompose into BCNF with all FDs preserved

```
algorithm Decomposition(set of elem. FDs. F, set of attributes A) : returns set {R<sub>i</sub>(A<sub>i</sub>, F<sub>i</sub>)}
 Result := \{R(A, F)\};
 Done := false;
Create F<sup>+</sup>;
 while not Done do
     if \exists R_i(F_i, A_i) \in \text{Result not being in BCNF then}
                                                                   // if there is a schema in the result violating BCNF
             Let X \rightarrow Y \in F_i such that X \rightarrow A_i \notin F^+.
                                                                  // X is not (super-)key and so X \rightarrow Y violates BCNF
                                                                  // we remove the schema being decomposed
             Result :=
                           (\text{Result} - \{\text{R}_{i}(\text{A}_{i}, \text{F}_{i})\}) \cup
                             \{R_i(A_i - Y, cover(F, A_i - Y))\} \cup // we add the schema being decomposed without attributes Y
                             \{R_i(X \cup Y, cover(F, X \cup Y))\} // we add the schema with attributes XY
     else
                                                               This partial decomposition on two tables is lossless, we get two schemas that
             Done := true;
                                                               both contain X, while the second one contains also Y and it holds X \rightarrow Y.
     endwhile
                                                               X is now in the second table a super-key and X \rightarrow Y is no more violating BCNF
 return Result;
                                                               (in the first table there is not Y anymore).
```

<u>Note</u>: Function cover(X, F) returns all FDs valid on attributes from X, i.e., a subset of F⁺ that contains only attributes from X. Therefore it is necessary to compute F⁺.

Example – decomposition

Contracts(A, F) A = {c = ContractId, s = SupplierId, j = ProjectId, d = DeptId, p = PartId, q = Quantity, v = Value} F = {c \rightarrow all, sd \rightarrow p, p \rightarrow d, jp \rightarrow c, j \rightarrow s}



The "Synthesis" algorithm

- algorithm for decomposition into 3NF, preserving dependencies
 - basic version not preserving lossless joins

algorithm **Synthesis**(set of elem. FDs F, set of attributes A) : **returns set** {R_i(F_i, A_i)} G = minimal cover of F compose FDs having equal left-hand side into a single FD every composed FD forms a scheme R_i (A_i, F_i) of decomposition **return** $\cup_{i=1..n}$ {R_i (A_i, F_i)}

- lossless joins can be preserved by adding another schema into the decomposition that contains *universal key*
 - i.e., a key from the original universal schema
- a schema in decomposition that is a subset of another one can be deleted
- we can try to merge schemas that have functionally equivalent keys, but such an operation can violate 3NF (or BCNF if achieved)!
 - i.e., we can try to minimize the number of relations

Example – synthesis

Contracts(A, F) A = {c = ContractId, s = SupplierId, j = ProjectId, d = DeptId, p = PartId, q = Quantity, v = Value} F = {c \rightarrow sjdpqv, sd \rightarrow p, p \rightarrow d, jp \rightarrow c, j \rightarrow s}

Minimal cover:

There are no redundant attributes in FDs.
Reundant FDs c → s and c → p were removed.
G = {c → j, c → d, c → q, c → v, sd → p, p → d, jp → c, j → s}

<u>Composition:</u> •G' = {c \rightarrow jdqv, sd \rightarrow p, p \rightarrow d, jp \rightarrow c, j \rightarrow s}

Result:

 $= \mathsf{R}_1(\{cqjdv\}, \{c \rightarrow jdqv\}), \ \mathsf{R}_2(\{sdp\}, \{sd \rightarrow p\}), \ \mathsf{R}_3(\{pd\}, \{p \rightarrow d\}), \ \mathsf{R}_4(\{jpc\}, \{jp \rightarrow c\}), \ \mathsf{R}_5(\{js\}, \{j \rightarrow s\}))$ (subset of R_2)

merging R_1 and R_4 (however, now $p \rightarrow d$ violates BCNF)