Lossless join



What does lossless join



usually, a schema is decomposed:

- when it is not in 3NF
- to improve efficiency
 - the smaller the size of the tuples, the greater the number of tuples we can load into memory in the same read operation
 - the tuple information is not used entirely by the operations on the database
 - example: the Student schema could be decomposed by separating the personal information (TaxC, Name, Surname, DateB, PlaceB, etc.) from the academic information (Matriculation, Degree, Year, etc.)

Lossless join



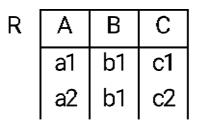
- •we have seen that when a schema is decomposed, **it is not enough for the** subschemas to be in 3NF
- ·let's see two examples (one more "abstract", the other more "concrete")



- •we now consider the schema R=ABC with the set of functional dependencies $F=\{A \rightarrow B, C \rightarrow B\}$ (the schema is not in 3NF due to the presence in F^+ of the partial dependencies $A \rightarrow B$ and $C \rightarrow B$, since the key is AC)
- •this schema can be decomposed into:
 - R1=AB with $A \rightarrow B$
 - R2=BC with $C \rightarrow B$
- •the resulting schemas, even if they **preserve all dependencies in** F^+ , is still not satisfactory.



•consider the **legal** instance of *R*:



the two facts (a1,b1,c1) and (a2,b1,c2) are true and not others

•we decompose it by obtaining:



•if we now compute the natural join we obtain:

| R | Α | В | С | |
|---|----|----|------------|--|
| | a1 | b1 | c 1 | |
| | a2 | b1 | c2 | |
| | a1 | b1 | c2 | these tuples did not exist in the original |
| | a2 | b1 | c1 | instance! |

 we must ensure that the decomposition and the following natural join does not result in any <u>loss of</u> <u>information</u>



- •consider the schema:
- •R={EmployeeID, ProjectID, Manager} with the set of functional dependencies :
- •F={EmployeeID → ProjectID, ProjectID → Manager
- •a project can have multiple managers but each manager has only one project, and an employee on a project reports to only one manager
- •the schema is not in 3NF due to the presence in F^+ of the partial dependencies $EmployeeID \rightarrow ProjectID$ and $ProjectID \rightarrow Manager$, since the key is (EmployeeID, Manager)
- •the schema can be decomposed into:
 - R1 = {Matriculation, Project} with EmployeeID → ProjectID and
 - R2 = {Project, Boss} with ProjectID → Manager
- •such a schema, while preserving all dependencies in F^+ , is not satisfactory



consider the legal instance of R:

| R | EmployeeID | ProjectID | Manager |
|---|------------|-----------|---------|
| | 501 | 30 | E1 |
| | 502 | 30 | E2 |

only the two facts (501,30,E1) and (501,30,E2) are true!

•based on the given decomposition, this instance decomposes into:

| D | 4 |
|-----------------------|---|
| $\boldsymbol{\Gamma}$ | |

| EmployeeID | ProjectID |
|------------|-----------|
| 501 | 30 |
| 502 | 30 |

| ProjectID | Manager |
|-----------|---------|
| 30 | E1 |
| 30 | E2 |



 ...and instead if you join the two legal instances resulting from the decomposition you get

| R | EmployeeID | ProjectID | Manager | |
|---|------------|-----------|---------|---|
| | 501 | 30 | E1 | |
| | 502 | 30 | E2 | |
| | 501 | 30 | E2 | <u> </u> |
| | 502 | 30 | E1 | tuples unrelated to the reality of interest |
| | | | | |
| | | | | |



in conclusion, the following requirement of the decomposed schema should be met:

 the decomposition must allow us to reconstruct, by natural join, each legal instance of the original schema (without adding any extraneous information)

Definition



if we decompose a relation schema R we want the obtained decomposition $\rho = R_1, R_2, ..., R_k$, such that each legal instance r of R can be reconstructed through natural join from the legal instances $r_1, r_2, ..., r_k$ of the decomposition schemas $R_1, R_2, ..., R_k$

as for reconstructing a tuple t of r it is required that $t[R_i] \in r_i$, $\forall i = 1,...,k$, then we must have $\pi_{R_i}(r) = r_i$, $\forall i = 1,...,k$

Definition

let R be a relation schema; A decomposition $\rho = R_1, R_2, ..., R_k$ has a **lossless join** if <u>for each legal instance</u> r of R we have

$$\mathbf{r} = \pi_{R1}(\mathbf{r}) \bowtie \pi_{R2}(\mathbf{r}) \bowtie ... \bowtie \pi_{Rk}(\mathbf{r})$$

Problem



 again, we start with a given decomposition, and look for a way to verify that it complains with the definition we gave

Theorem



Theorem: let R be a relation scheme and let $\rho = R_1, R_2, ..., R_k$ be a decomposition of R; for each legal instance r of R, denoted by $m_{\rho}(\rho) = \pi_{R1}(r) \bowtie \pi_{R2}(r) \bowtie ... \bowtie \pi_{Rk}(r)$, we have:

- \bullet r \subseteq $m_{\rho}(r)$
- $\bullet \pi_{Ri}(m_{\rho}(r)) = \pi_{Ri}(r)$
- $\bullet m_{\rho}(m_{\rho}(r)) = m_{\rho}(r)$

Checking lossless join



- given a relation scheme R, a set of functional dependencies F, and a decomposition ρ
- how do we check that the given decomposition has a lossless join?
- there exist an algorithm that checks this in polynomial time
- we need F as we are implicitly checking legal instances of R

Algorithm



Algorithm - lossless join

Input: a relation scheme **R**, a set **F** of functional dependencies on R, a decomposition

 $\rho = R_1, R_2, ..., R_k$

Output: whether or not ρ has a lossless join

begin

construct a table r as follows:

r has R columns and k rows

at the intersection of the i-th row and j-th column put:

the symbol a_j if the attribute $A \in R_{i,j}$

the symbol b_{ii} otherwise

repeat

for every $X \rightarrow Y \in F$

the attribute A, is part of the subschema R,

a column for each attribute of R and a row for each

element of the decomposition (subschema)

Index i = element of the decomposition = line

Index j = attribute = column

if there are two tuples t_1 and t_2 in r such that $t_1[X] = t_2[X]$ and $t_1[Y] \neq t_2[Y]$ then for every attribute A_i in Y: if $t_1[A_i] = a_i$ then $t_2[A_i] = t_1[A_i]$

else $t_1[A_i] = t_2[A_i]$

until r has a line with all "a" or r has not changed;

if r has a row with all "a" then: it has a lossless join

else: it does not have a lossless join

even in this case the algorithm <u>always</u> ends! then we need to check if in r there is the tuple we are looking for

we also correctly handle the

case where $t_{2}[A_{1}]='a_{1}$

end

Comments



- we can consider the a_j as particular values belonging to the domain of the attribute A_j
- we can consider $\mathbf{b_{i,j}}$ as particular values belonging to the domain of the attribute $\mathbf{A_j}$
- we can consider all values a equal to each other
- the value $\mathbf{b_{ij}}$ is different from $\mathbf{a_j}$ and from another value $\mathbf{b_{kj}}$, even if they all belong to the same domain (that of attribute $\mathbf{A_i}$)
- as a consequence, the initial r is a particular instance
 of the schema R

Comments



- the algorithm modifies r, so that all dependencies in F are fulfilled
- whenever it finds two tuples that are equal on the determinant but different on the dependent, it modifies the dependent so they become equal
- in doing so, it prioritizes the "a" symbol ("a" never becomes "b"; "b" can become "a")

Comments



- if two tuples have the same value in the determinant but different values in the dependent of a dependency, and only one of the tuples has the value "a" in the dependent, we change the "b" of the other tuple into an "a"
- if two tuples have the same value in the determinant but different values in the dependent of a dependency, and neither has a value of "a" in the dependent, we change one of the tuples so that they have the same value of "b" (e.g., if we have b_{ij} and b_{kj} then we make both values b_{ij} or b_{kj})
- two values are equal if they are both "a" or if they have a "b" with the same subscript
- the algorithm stops when **ALL** pairs of tuples satisfy all the dependencies in F
- so, in the end, r became a legal instance of R

Theorem



Theorem: let R be a relation scheme, F a set of functional dependencies on R and let $\rho = R_1, R_2, ..., R_k$ be a decomposition of R; the algorithm correctly decides whether ρ has a lossless join

- Demonstration: it must be shown that:
 - has a lossless join $(m_{\rho}(r)=r)$, for each legal instance r)
 - if and only if
- when the algorithm terminates, r has a tuple with all "a"

Theorem: proof



if (by contradiction)

•suppose **that** ρ **has a** lossless **join** $(m_{\rho}(r)=r)$ and that when the algorithm terminates the table r does not have any tuple with all "a"

*r is a legal instance of R, since the algorithm terminates when there are no more dependency violations in F

•as no "a" symbol appearing in the initial value of r was changed into a "b" by the algorithm, for each i, i=1,...,k, $\pi_{Ri}(r)$ contains (from the beginning!) a tuple with all 'a', which was the one obtained by projecting the instance r on the attributes of R_i , more precisely in the row corresponding to the subschema R_i so, $m_{\rho}(r)$ contains a tuple with all 'a' and, consequently, $m_{\rho}(r) \neq r$

rso, *m_ρ(r) contains a tuple with all a* and, consequently, *m_ρ(r)≠r* (contradiction)