

Preserving F



SAPIENZA
UNIVERSITÀ DI ROMA

What it means to preserving dependencies



- a schema is usually decomposed:
- **when it is not in 3NF**
- for efficiency reasons:
 - the smaller the size of the tuples, the greater the number we can load into memory in the same read operation
 - if the tuple information is not used by the same type of operations in the database, it is better to decompose the schema
 - Example: the Student schema could be decomposed by separating the personal information (CF, FirstName, LastName, DateBirth, PlaceBirth, etc.) from the academic information (Matriculation, CourseDegree, YearCourse, etc.).

What it means to preserving dependencies



SAPIENZA
UNIVERSITÀ DI ROMA
DIPARTIMENTO DI INFORMATICA

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

What we want to achieve - 3NF is not enough



- a schema that is not in 3NF can be **decomposed** in **multiple ways** into a set of schemas in 3NF
- for example, the schema $R = ABC$ with the set of functional dependencies $F\{A \rightarrow B, B \rightarrow C\}$ is not in 3NF, due to the presence in F^+ of the transitive dependency $B \rightarrow C$ (the key is A)
- R can be decomposed into:
 - $R1=AB$ with $A \rightarrow B$ and
 - $R2=BC$ with $B \rightarrow C$
 - **or**
 - $R1=AB$ with $A \rightarrow B$ and
 - $R2=AC$ with $A \rightarrow C$
- both schemas **are** in 3NF, however the second solution **is not satisfactory**

What we want to achieve - 3NF is not enough



- if we consider two **legal** instances of the obtained schemas:

R1	A	B
	a1	b1
	a2	b1

R2	A	C
	a1	c1
	a2	c2

- the instance of the original R schema that I can reconstruct from this (the only way is to reconstruct it by doing a natural join!) is:

R	A	B	C
	a1	b1	c1
	a2	b1	c2

- BUT** it is not a legal instance of R , since it **does not satisfy** the functional dependency $B \rightarrow C$

we want to preserve ALL dependencies in F^+

Example



- let us consider the schema $R = \{Matriculation, Town, Province\}$ with the set of functional dependencies:
- $F = \{Matriculation \rightarrow Town, Town \rightarrow Province\}$
- the schema is not in 3NF due to the presence in F^+ of the transitive dependency $Town \rightarrow Province$ (the key is $Matriculation$ and $Province$ transitively depends on $Matriculation$)
- R can be decomposed into:
 - $R1(Matriculation, Town)$ with $Matriculation \rightarrow Town$
 - $R2(Town, Province)$ with $Town \rightarrow Province$
- **or**
- $R1(Matriculation, Town)$ with $Matriculation \rightarrow Town$
- $R2(Matriculation, Province)$ with $Matriculation \rightarrow Province$
- both schemas **are** in 3NF, but the second solution **is not satisfactory**

What we want to achieve - 3NF is not enough



- consider the **legal** instances of the obtained schemas:

R1	Matriculation	Town
	501	Tivoli
	502	Tivoli

R2	Matriculation	Province
	501	Rome
	502	Rieti

- the instance of the original R schema that we can reconstruct from this with a natural join is:

R	Matriculation	Town	Province
	501	Tivoli	Rome
	502	Tivoli	Rieti

- but** it is not a legal instance of R , since it **does not satisfy** the functional dependency $Town \rightarrow Province$
- clearly there was an error in the data, but we could not detect it**

- **Definition:** let R be a relation scheme; a decomposition of R is a family R_1, R_2, \dots, R_k of subsets of R **covering** R , that is

$$R = \bigcup_{i=1}^k R_i$$

(the subsets may have **nonempty** intersection)

- **in other words:** decomposing it means defining **subschemas**, **each of them** containing **a subset of** the attributes of R and their union is R
- **so,** R is a set of attributes and a decomposition of R is a family of sets of attributes

Example



- given the schema $R = (\text{TaxC}, \text{Name}, \text{Surname}, \text{Matriculation}, \text{DateB}, \text{PlaceB}, \text{Degree}, \text{Year})$ we can have the decompositions:
 - $R_1 = (\text{TaxC}, \text{Matriculation}, \text{Name}, \text{Surname}, \text{DateB}, \text{PlaceB})$
 - $R_2 = (\text{Matriculation}, \text{TaxC}, \text{Name}, \text{Surname}, \text{Degree}, \text{Year})$
 - or
 - $R_1 = (\text{TaxC}, \text{Matriculation}, \text{Name}, \text{Surname}, \text{DateB}, \text{PlaceB})$
 - $R_2 = (\text{Matriculation}, \text{Degree}, \text{Year})$
- in both cases, the families of subsets of R (subschemas) cover the schema, but then we need to check if both decompositions are "good", i.e., if the subschemas are in 3NF, F is preserved, and we can reconstruct any legal instance of R by natural join of the instances of the subschemas



- **"decomposing" an instance of a relation with a certain schema, based on the decomposition of the schema itself, means projecting each tuple of the original instance onto the attributes of the individual subschemes ...**
- **...eliminating the duplicates** that might be generated by the fact that two tuples are distinct **but have a common portion** that falls into the same subschema

Example



R

TaxC	Name	Last name	Matriculation	Datan	PlaceB	Degree	Year
DDD	Davide	Bigi	1111	12-12-90	Bari	Letters	3
FFF	Gianni	Neri	1212	15-09-91	Milan	Law	2
AAA	Antonio	Rossi	1313	09-08-92	Naples	Mathematics	1

R

1

TaxC	Name	Surname	Matriculation	DateB	PlaceB
DDD	Davide	Bigi	1111	12-12-90	Bari
FFF	Gianni	Neri	1212	15-09-91	Milan
AAA	Antonio	Rossi	1313	09-08-92	Naples

subschema of R

projection of the instance
onto the subschema

R

2

TaxC	Name	Surname	Matriculation	Degree	Year
DDD	Davide	Bigi	1111	Letters	3
FFF	Gianni	Neri	1212	Law	2
AAA	Antonio	Rossi	1313	Mathematics	1

subschema of R

projection of the instance
onto the subschema

.we continue with the definition of **equivalence between two sets of functional dependencies**.

.Definition: let F and G be two sets of functional dependencies;
 F and G are **equivalent**, written $F \equiv G$ if

$$F^+ = G^+$$

.F and G do NOT contain the same dependencies, but their closures do

• **checking the equivalence of** two sets of functional dependencies F and G requires verifying the **equality of F^+ and G^+** , i.e., that

$$F^+ \subseteq G^+ \text{ and that } F^+ \supseteq G^+$$

- as mentioned earlier, computing the closure of a set of functional dependencies takes **exponential time**
- however, the following **lemma** and an algorithm allow us to check the equivalence of the two sets of functional dependencies **in polynomial time**

- **Lemma:** let F and G be two sets of functional dependencies;
if $F \subseteq G^+$ then $F^+ \subseteq G^+$
- **Demonstration:** let $f \in F^+ - F$ (f is a dependency in F^+ that **does not appear** in F)
- **every functional dependency in F can be derived from G** by Armstrong's axioms (for **the hypothesis** it is in G^+ , and we know that $F^+ = F^A$)
- always for the Theorem that proves that $F^+ = F^A$, $f \in F^+$ can be derived **from F by applying the Armstrong's axioms**
- **so, any $f \in F^+$ can be derived from G by Armstrong's axioms**

$$G \xrightarrow{A} F \xrightarrow{A} F^+$$

- **Definition:** let **R** be a **relation scheme**, **F** a **set of functional dependencies on R** and $\rho = R_1, R_2, \dots, R_k$ a **decomposition** of R:
- we say that ρ **preserves F** if

$$F \equiv G = \bigcup_{i=1}^k \pi_{R_i}(F),$$

where $\pi_{R_i}(F) = \{ X \rightarrow Y \in F^+ \mid XY \subseteq R_i \}$

- **note:**
- obviously, $G = \bigcup_{i=1}^k \pi_{R_i}(F)$ is a set of functional dependencies
- each $\pi_{R_i}(F)$ is a set of functional dependencies obtained by projecting F onto the subschema R_i
- projecting a set of dependencies F onto a subschema R_i means **taking all** the dependencies **derivable from F** by the Armstrong's axioms (so, those **in F⁺**) that have **all** the attributes (determinants and dependents) **in R_i**

Check if F is preserved



- suppose we **already have** a **decomposition** and want to **check if it preserves the functional dependencies** in F
- checking **whether** a decomposition preserves a set of functional dependencies F requires verifying **the equivalence** of the two sets of functional dependencies F and $G = \bigcup_{i=1}^k \pi_{R_i}(F)$
- **note: because the definition of G**, it will always be that $F^+ \supseteq G$
- each **projection of F** that is **included in G** is also in F^+ , and, by the lemma, that implies that $G^+ \subseteq F^+$
- so, given the Lemma, it is enough to verify that $F \subseteq G^+$

Check if F is preserved



Algorithm - F in G^+

Input: two sets F and G of functional dependencies on R ;

Output: the success variable (**true** if $F \subseteq G^+$, **false** otherwise)

begin

success = true

for each $X \rightarrow Y \in F$

begin

compute X_G^+

if $Y \notin X_G^+$ **then** *success = false*

end

end

- **if** $Y \subseteq X_G^+$ **then** $X \rightarrow Y \in G^A$ (for the lemma) and $X \rightarrow Y \in G^+$ (for the Theorem)
- it is enough to verify that **even a single dependency** does not belong to the closure of G to be able to state that **the equivalence does not exist**



- **how do we compute X_G^+ ?**
- to use the Algorithm already seen for calculating the closure of X , we would **first need to calculate G**
- ... but, **by the definition of G** , this **requires** the computation of F^+ , which takes **exponential** time
- we present an algorithm to compute X_G^+ **from F**

Calculating X_G^+ from F



Algorithm - X_G^+ from F

Input: a schema R , a set F of functional dependencies on R , a decomposition $\rho = R_1, R_2, \dots, R_k$, a subset X of R

Output: the closure of X with respect to $G = \bigcup_{i=1}^k \pi_{R_i}(F)$, (in variable Z)

begin

$Z = X$

$S = \emptyset$

for $i = 1$ **to** k

$S = S \cup (Z \cap R_i)^+_F \cap R_i$

while $S \not\subseteq Z$

begin

$Z = Z \cup S$

for $i = 1$ **to** k

$S = S \cup (Z \cap R_i)^+_F \cap R_i$

end

end

we note the relationship between the definition of the projections of F :

$$\pi_{R_i}(F) = \{ X \rightarrow Y \in F^+ \mid XY \subseteq R_i \}$$

and the update step of S :

$$S = S \cup (Z \cap R_i)^+_F \cap R_i$$

we are collecting the attributes functionally determined **by the part of Z** that intersects R_i , (for the lemma, $A \in X^+$ if and only if $X \rightarrow A \in F^+$, **and in this case $X = Z \subseteq R_i$**) according to the dependencies in F^+ ; and from these we select those contained in R_i ,

starting from a set X , we are "reconstructing" its **closure** from the projections of F that **make up G** , and from there iterating from the set G^+ (through transitivity)

by intersecting with R_i , we ensure that the dependency is valid in the individual subschema

- we also note the difference with the update step of S in the two algorithms:

$$S = A \mid Y \rightarrow V \in F, A \in V \wedge Y \subseteq Z$$

$$S = S \cup (Z \cap R_i)^+_{F_i} \cap R_i$$

- here the union is needed as **we are considering** the different subschemas **in turn**
- we also take dependencies in the closure **$F^+ (Z \cap R_i \rightarrow A \in F^+$ and... the lemma)** as **the dependencies in the set G** are included in the projections of F onto the subschemas, which are however, dependencies in **F^+**
- in conclusion we will find the attributes that functionally depend on X even if they belong to subschemas in which X is not included! as they depend on attributes that are in the same subschema as X and depend on X, but they are also in other subschemas!

- the algorithm (by definition) **always terminates**
- the fact that the algorithm terminates **does not indicate** that a dependency $X \rightarrow Y$ **is preserved!**
- to **check if** $X \rightarrow Y$ is preserved, we need to use the previous algorithm

- **Theorem:** let R be a relation schema, F a set of functional dependencies on R , $\rho = R_1, R_2, \dots, R_k$ a decomposition of R and X a subset of R ; the Algorithm correctly computes X_G^+ , where $G = \pi_{R_i}(F)$