Normalization Thanks to Lois Delcambre

Normalization Strand Map

© Lois Delcambre, David Maier 2005-2018

All computer science students must learn to integrate theory and practice, to recognize the importance of abstraction, and to appreciate the value of good engineering design. cc2001 (http://www.sigcse.org/cc2001/)

- This normalization strand map has two strands:
 - the practical concepts and techniques and
 - the formal concepts and results for normalization.



3

Normalization based on Functional Dependencies (FDs)



2

- *Normalization* works on relations in a relational database.
- *Normalization* based on FDs has been formalized and the theory has been completely worked out.
- Normalization makes it easier to update a database but may degrade query performance – an engineering tradeoff.

© Lois Delcambre, David Maier 2005-2018



© Lois Delcambre, David Maier 2005-2018



Functional Dependencies (FDs) generalize keys

Functional dependencies (FDs) for relational tables are a generalization of the notion of key for a table.



7

Functional Dependencies

8

- An FD, X \rightarrow Y, where X is a set of attributes and Y is a set of attributes
- It is a statement that if two tuples agree on attributes X they must agree on attributes Y

For each X value there is ONLY one Y value.

Functional Dependencies (from semantics of the application)



Likely functional dependencies: $ssn \rightarrow employee$ -name course-number dept \rightarrow course-title

Unlikely functional dependencies dept ¥ book birthdate 🔆 ssn

Will FDs be enforced?



10

12

Consider this table: Emp(ssn, name, phone, dnum, dept-name)

Suppose there is an FD from $dnum \rightarrow dept-name$

But ssn is the key for this table. What will prevent two names for one dept?

© Lois Delcambre, David Maier 2005-2018

Will this FD be enforced? Let's try it.

Consider this table: Emp(ssn, name, phone, dnum, dept-name)

Employee	ssn	Nam	Phone	Dnum	Dept name
		е			$r \rightarrow$
	111111111	John	555-1234	12	Sales
	222222222	Mary	555-7890	12	Marketing

Can we put these two rows in this table?

Yes, it doesn't violate the key constraint.

But, the FD from dept to dept-name is violated! We shouldn't haven't two different names for dnum 12!



11

9

Functional Dependencies

© Lois Delcambre, David Maier 2005-2018

For an FD

 $X \rightarrow Y$

We say that X determines Y

We want to know if it is always true in the application.

We can then use FDs to figure out the keys for tables and to normalize the tables.

```
© Lois Delcambre, David Maier 2005-2018
```

Every key implies a set of FDs



13

Each key implies a set of functional dependencies (FDs) from the key to the non-key attributes.

Employee (ssn, name, salary, job-code)

FDs implied by the key: $ssn \rightarrow name$ $ssn \rightarrow salary$ $ssn \rightarrow job-code$

© Lois Delcambre, David Maier 2005-2018

The Problem: "Troublesome" FDs

"Troublesome" FDs (FD where the left-hand-side of the FD is NOT a key for the table where its attributes appear) cause redundancy and update anomalies.



15

© Lois Delcambre, David Maler 2005-2018

But, some FDs are NOT implied by the key.





© Lois Delcambre, David Maier 2005-2018

Advantages & Disadvantages of Redundancy



16

- Disadvantage: Any time information is stored more than once, it has the possibility of being inconsistent.
 - Phone numbers in your laptop
 - Phone numbers in your cell phone
 - Phone numbers in your address book
- If someone changes his or her phone number, do you remember to change it in every place?
- Advantage: Redundant information may improve retrieval speed

Sometimes Redundancy is Caused by FDs

© Lois Delcambre, David Maier 2005-2018



17

19

Consider this table:



Then *dname* and *dmgr* are stored redundantly – whenever there are multiple employees in a department.

This redundancy is caused by what we informally call "troublesome" FDs. The FDs shown in blue are "troublesome".

Redundancy Caused by Troublesome FD – Sample Data

•••

. . . .

EMP(name	, <u>ssn</u> , l	oirthdate,	address,	dnum, o	dname, o	dmgr)
John	111	June 3	123 St.	D1	sales	222
Sue	222	May 15	455 St.	D1	sales	222
Max	333	Mar. 5	678 St.	D2	researc	ch 333
Wei	444	May 2	999 St.	D2	researc	h 333
Tom	555	June 22	888 St.	D2	researd	h 333

We have the department name and manager twice for D1 and three times for D2!



EMP(name, ssn, birthdate, address, dnum, dname, dmgr)

The name and the manager of the department is repeated, for each employee that works in that department.

Redundancy!

If you replicate information, the copies might be inconsistent.

Update Anomalies

caused by "troublesome" FDs EMP(name, <u>ssn</u>, birthdate, address, dnum, dname, dmgr)

© Lois Delcambre, David Maier 2005-2018

r)

20

Insertion anomalies:

if you insert an employee with a department then you need to know the descriptive information for that department.

if you want to insert a department, you can't ... until there is at least one employee.

Deletion anomalies: if you delete an employee, is that dept. gone? Was this the last employee in that dept.?

Modification anomalies: If you want to change *dname*, for example, you need to change it everywhere! And you have to find them all first.

Troublesome FDs cause (redundancy and) update anomalies.

© Lois Delcambre, David Maier 2005-2018

Aside: Null Values – Advantages & Disadvantages

Null values can help with some update issues but they can make SQL queries more complex.



Null Values are Useful

Consider



Employee(<u>ssn</u>, name, DOB, partner)

Null values make it simpler to insert data, for example

• before you know the partner name しんいん

• when there is no partner

Allowing null values makes the DB more flexible.

DNF

does not exist

© Lois Delcambre, David Maier 2005-2018

Null Values Cause Problems for Aggregate Operators



21

Employee (<u>ssn</u>, name, salary)

SELECT AVG(salary) FROM Employee;

SELECT SUM(salary) INTO salsum FROM Employee; SELECT COUNT(*) INTO total FROM Employee;

salsum/total might be different from first query answer. How could that happen?

They can also complicate joins and WHERE clauses

© Lois Delcambre, David Maier 2005-2018

23

Splitting Tables Reduces the Need for Null Values

© Lois Delcambre, David Maier 2005-2018

24

22

Use two tables: Employee (<u>ssn</u>, name, DOB) Employee-extra (<u>ssn</u>, partner)

Rather than: Employee(<u>ssn</u>, name, DOB, partner)

© Lois Delcambre, David Maier 2005-2018

Generally, it is better to reduce the use of null values, if you can. The first design, above, doesn't require the use of null values for partner.

The Solution: Lifting "Troublesome" FDs

Normalization by decomposition, based on FDs (where "troublesome" FDs are lifted into a separate table), reduces redundancy and update anomalies.



25

© Lois Delcambre, David Maier 2005-2018



Example: Finding Troublesome FDs

26

28

Table is Split onto New Schemas

© Lois Delcambre, David Maier 2005-2018

New-EMP(name, ssn, birthdate, address, dnum)

`	·	,	,	,
John	111	June 3	123 St.	D1
Sue	222	May 15	455 St.	D1
Max	333	Mar. 5	678 St.	D2
Wei	444	May 2	999 St.	D2
Tom	555	June 22	888 St.	D2

Dept(dnum, dname, dmgr)

anomalies

D1	sales	222
D2	research	333

Less redundancy! Tastes better! Fewer update issues!

Basic Idea: Normalize based on FDs

- Identify all the (non-trivial) FDs in an application.
- Identify FDs that are implied by the keys.
- Identify FDs that are NOT implied by the keys the "troublesome" ones.
- Decompose a table with a "troublesome" FD into two or more tables by "lifting" each troublesome FD into a table of its own. Note: when there are two or more "troublesome" FDs with the same left side, then they can be lifted, together, into a single table.

29

31

© Lois Delcambre, David Maier 2005-2018

Advantages of Normalization based on Decomposition

When this table:

Emp(name, ssn, birthdate, address, dnum, dname, dmgr)

is replaced by these two tables:

© Lois Delcambre, David Maier 2005-2018

Dept(d<u>num</u>, dname, dmgr) New-Emp(name, <u>ssn</u>, birthdate, address, dnum)

Are there any update anomalies in the new tables?





Emp(name, <u>ssn</u>, birthdate, address, dnum, dname, dmgr) split into

Dept(d<u>num</u>, dname, dmgr)

New-Emp(name, ssn, birthdate, address, dnum)

If there are applications that currently query Emp, can define a view: CREATE VIEW Emp AS SELECT *

FROM Dept NATURAL JOIN New-Emp

© Lois Delcambre, David Maier 2005-2018

Update statements will require changes in most cases

Let's Check the Update Anomalies

•••
• •

30

Insertion anomalies:

if you insert an employee with a department then you need to know the descriptive information for that department. NO – ONLY THE NUMBER if you want to insert a department, you can't ... until there is at least one employee. NO PROBLEM

Deletion anomalies: if you delete an employee, is that dept. gone? Was this the last employee in that dept.? NO PROBLEM

Modification anomalies: If you want to change *dname*, for example, you need to change it everywhere! And you have to find them all first. dname is only stored once! Is there any redundancy? Some – in the foreign key.

© Lois Delcambre, David Maier 2005-2018

Questions about normalization

- How do we know which FDs we have? Talk to domain experts; identify FDs; use them as the starting point for normalization.
- How do we know if the decomposition is correct?
- How do we know how much to normalize? How far should we go?
- How do we know if all of the FDs of interest are being enforced – by using keys for tables?
 We need the formal definition of FDs to be able to answer these questions.

© Lois Delcambre, David Maier 2005-2018

Definition of a function

Remember the definition of a function:

X	f(x)	х	g(x)	х	h(x)
(1)	2	1	2	1	10
1	3 /	2	2	2	20
2	5	3	5	3	30
3	5				

Which of these are functions?

An FD is a functional relationship (that always holds in a relation) among attribute values

© Lois Delcambre, David Maier 2005-2018



35

33

FDs and Keys: Formal Definition

A functional dependency is formally defined as a functional relationship between two sets of attributes. This leads to the definition of trivial FDs and superkeys.



© Lois Delcambre, David Maier 2005-2018

34

Answer (concerning functions)



36

Remember the definition of a function:

Х	f(x)	х	g(x)	х	h(x)
1	2	1	2	1	10
1	3	2	2	2	20
2	5	3	5	3	30
3	5				

f is NOT a function because for an input of "1" there are two answers ("2" and "3"). g and h are functions.

Example of an FD – a function



37

39

Employee (ssn, name, phone, salary)

Since $ssn \rightarrow name$ is an FD If we know that there is only one name for an ssn,

then we know that $ssn \rightarrow name$ is a function

We don't expect *salary* \rightarrow *phone* to be function.

© Lois Delcambre, David Maier 2005-2018

Trivial FD

We have a trivial FD whenever the attributes on the right side of an FD are a subset of the attributes on the left side of the FD:

name phone \rightarrow phone

Trivial FDs aren't "troublesome" and won't help us decompose a table. *Ignore them.*

Another Example



Employee (ssn, name, phone, dept, dept-mgr)

$dept \rightarrow dept-mgr$

If we know that there is only one dept-mgr for a dept,

then we know that $dept \rightarrow dept-mgr$ is a function!

© Lois Delcambre, David Maier 2005-2018

Definition of a Superkey for a Relation

•••
$\bullet \bullet \bullet \bullet$
• •

38

A *superkey* is a set of attributes from a relation that contains a key.

Every key is (automatically) a superkey. A superkey is NOT necessarily a key.

Example:

Emp (<u>ssn</u>, name, phone, dept) ssn is a key (and hence a superkey) for this relation. (dept, ssn) is a superkey for this relation (but not a key).

© Lois Delcambre, David Maier 2005-2018

© Lois Delcambre, David Maier 2005-2018



A key is a minimal set of attributes in a relation whose values are guaranteed to uniquely identify tuples in the relation.

- Two distinct tuples have distinct key values
- (minimal) No subset of the fields that comprise a key is a key
- Can be more than one key for a table (not just a single declared key as in SQL)

41

43

Keys and FDs are Constraints



- We need to know if keys and FDs (always) hold in the application.
- We need to consult a domain expert to find out what the keys and FDs are. The keys and FDs serve as input to the database design process.

2NF, 3NF, BCNF: Normal forms based on FDs

© Lois Delcambre, David Maier 2005-2018

Given a set of FDs and one or more tables, three increasingly stronger normal forms, namely 2NF, 3NF, and BCNF, have been defined. BCNF implies 3NF.

3NF implies 2NF.

(BCNF = Boyce-Codd Normal Form)

© Lois Delcambre, David Maier 2005-2018

Informal Definitions Normal Forms Based on FDs

© Lois Delcambre, David Maier 2005-2018



42

1NF - all attribute values (domain values) are atomic
(part of the definition of the relational model)2NF - all non-key* attributes must depend on a whole key (no
partial dependencies)
r ($\underline{A} \ \underline{B} \ \underline{C} \ \underline{D} \ \underline{E}$) $\underline{B} \rightarrow \underline{C}$ violates 2NF3NF - table is in 2NF and all non-key attributes must depend
on only a key (no transitive dependencies)
r ($\underline{A} \ \underline{B} \ \underline{C} \ \underline{D} \ \underline{E}$) $\underline{C} \rightarrow \underline{D}$ violates 3NFBCNF - every left side of an FD is a key for the table
(All FDs are implied by the keys)
r ($\underline{A} \ \underline{B} \ \underline{C} \ \underline{D} \ \underline{E}$) $\underline{C} \rightarrow A$ violates BCNF (but not 3NF)* "non-key" = not part of any key

Examples of Violations

2NF - all non-key attributes must depend on a whole key Assigned-to (<u>a-project, a-emp</u>, emp-name, percent)

3NF - all non-key attributes must depend on only a key Employee (<u>ssn</u>, name, address, project, p-title)

BCNF - every determinant (LHS of an FD) is a key for this table (all FDs are implied by the keys) emp-ID \rightarrow ssn

Assigned-to (emp-ID, a-project, ssn, percent)



© Lois Delcambre, David Maier 2005-201

Formal definition of BCNF



47

45

For a table R, every FD X \rightarrow A that occurs among attributes of R then either:

- A is an element of X (X \rightarrow A is trivial), or
- X is a superkey of R

For 3NF there is one other option:

© Lois Delcambre, David Maier 2005-2018

A is part of a key





Assigned-to (<u>a-project a-emp</u> emp-name, percent)

Employee (a-emp, emp-name)

- 1. Lift the troublesome FD(s) into a table of its own. Key for new table is left hand side of the troublesome FD.
- 2. Leave the left side of the FD behind in the original table. Assigned-to (<u>a-project (a-emp)</u> percent)
- 3. Eliminate emp-name from the Assigned-to table.

© Lois Delcambre, David Maier 2005-2018

Dependency Preservation: Using a sound & complete set of inference rules

Dependency preservation requires the use of a sound and complete set of rules of inference to compute F⁺, the *closure* of a set F of FDs. Given the original set of FDs, F. Let G consist of the FDs in F+ whose attributes appear in any relation scheme (after normalization). Dependency preservation is when F⁺ = G⁺.



48

Example

Consider the following table: Employee (ssn, name, phone, dept, dept-name)

Original FDs F: ssn→name ssn→phone ssn→dept ssn→dept-name dept→dept-name Employee (ssn, name, phone, dept) Department (dept, dname) drof - name Resulting FDs G:

 $ssn \rightarrow name ssn \rightarrow phone$

dept→dept-name

What about ssn→dept-name? Is it lost? Can there be two department names for one ssn?

ssn→dept

© Lois Delcambre, David Maier 2005-2018

We need to derive all FDs from a given set of FDs. We need rules. sets of attributes



51

Complete

49

For sets of attribute X and Y Reflexivity

If Y is a subset of X, then $X \to Y$

As an example, for all attributes, $A \rightarrow A$

union examples: name \rightarrow name, gender \rightarrow gender

Augmentation

If FD X \rightarrow Y holds, then so does XZ $\xrightarrow{}$ YZ, for all Z As an example, augmentation creates superkeys from keys. Sound and Transitivity

If FDs $X \rightarrow Y$ and $Y \rightarrow Z$ hold. then so does $X \rightarrow Z$

© Lois Delcambre, David Maier 2005-2018





50

52

Employee (ssn, name, phone, dept, dept-name) Original FDs F: ssn→name ssn→phone ssn→dept dept→dept-name ssn→dept-name Employee (ssn, name, phone, dept) Department (dept, dname) Resulting FDs G: ssn→name ssn→phone ssn→dept dept→dept-name What about ssn-dept-name? Is it lost? Can there be two department names for one ssn?

© Lois Delcambre, David Maier 2005-2018

NO! It's not lost. One ssn has only one dept. And one dept has only one dept-name. So ssn has only one dept-name.

Use the rules to compute closure

Let F be a set of FDs.

F⁺ is the set of all FDs implied (or derivable) from F using a sound & complete set of inference rules **Reflexivity**: If Y is a set of attrs, Y subset of X, then $X \rightarrow Y$ Augmentation: If $X \rightarrow Y$, and Z is a set of attrs, then $XZ \rightarrow XY$ **Transitivity**: If $X \rightarrow Y$, $Y \rightarrow Z$, then $X \rightarrow Z$

Compute F⁺ by applying rules until no new FDs arise. F^+ is called the closure of F_{h}^+

Definition of Dependency Preserving Deropole (R)



Suppose F is the original set of FDs.

Compute F+.

G is set of FDs from F⁺ that are present in individual relations in \emptyset . D

Compute G⁺.

If $F^+ = G^+$

then the decomposition is dependency preserving

For a complex design, you may want to implement one of the known algorithms for testing $F^+ = G^+$.

© Lois Delcambre, David Maier 2005-2018

Decompose: Project Operator "Recompose": Join Operator

When

Emp(name, ssn, birthdate, address, dnum, dname, dmgr)

is replaced by these two tables: Department(d<u>num</u>, dname, dmgr) NewEmp (name, <u>ssn</u>, birthdate, address, dnum)

We use the project operator to decompose

Department = $\pi_{dnum,dname,dmgr}$ Emp

NewEmp = $\pi_{name,ssn,birthdate,address,dnum}$ Emp And we use the join operator to put the pieces together Emp = Department $\bowtie_{D.dnum=NE.dnum}$ NewEmp



53

Decomposition is correct when it is lossless

The decomposition algorithm (based on lifting "troublesome" FDs into a separate table) guarantees that the decomposition of the original table is lossless.



© Lois Delcambre, David Maier 2005-2018

54

What is a lossless (and a lossy) decomposition?

56

We want to make sure that we haven't thrown away any information from the original schema.

When table R is decomposed into tables R1 and R2 then the decomposition is lossless (correct) if:

 $(R1 \bowtie R2)$ is identical to R natural join

If it is a lossy decomposition, then R1 ⋈ R2 gives you TOO MANY tuples.

© Lois Delcambre, David Maier 2005-2018



Is the Decomposition Algorithm Lossless?



- 1. Lift the "troublesome" FD(s) (all the FDs with the same LHS) into a table of their own. Key for new table is left hand side of the troublesome FD(s).
- 2. Leave the left side of the FD behind in the original table.
- 3. Eliminate the RHS attributes from the original table.
- Yes, we are guaranteed that the decomposition is lossless. The attribute in common is definitely a key for the new "lifted" table.

One Guarantee for a Lossless Decomposition



Consider a table: R (\underline{a} , b, c, d, e) with a troublesome FD d \rightarrow e. Decompose it into two tables: R₁(\underline{a} , b, c, d) R₂(\underline{d} , e)

As long as

the attributes in common are a key for (at least) one of the relations, R_1 or R_2 then we know that the decomposition is lossless. For this example d is the attribute in common. And d is a key for R_2 , the second table.

Example of a Lossy Decomposition (revisited)

© Lois Delcambre, David Maier 2005-2018

© Lois Delcambre, David Maier 2005-2018

60

58

Employee(<u>emp-number</u>, name, p-num, p-title)

decomposition: Employee (<u>emp-number</u>, name) Project (<u>p-num</u>, p-title, <u>name</u>)

Notice that the common attribute, name, is not a key for either of these tables.





lossless decomposition

don't throw any information away be able to reconstruct the original relation

dependency preservation

all of the original, non-trivial FDs can be derived from FDs implied by the keys of resulting tables

Boyce-Codd normal form (BCNF) - no redundancy beyond foreign keys; all FDs implied by keys

© Lois Delcambre, David Maier 2005-2018

Counterexample (a table that can't be decomposed into BCNF with dependency preservation)

63

61

Original table - a table that holds US addresses

addr(number street city state zip)

© Lois Delcambre, David Maier 2005-2018

The original FDs are: number street city state \rightarrow zip zip \rightarrow state



It is not always possible to have BCNF AND dependency preservation

62

64

© Lois Delcambre, David Maier 2005-2018





Let's decompose addr(<u>number street city state</u> zip) using this "troublesome" FD:

 $zip \rightarrow state$

Addr2 (<u>number, street, city, zip</u>) Zip-state (<u>zip</u>, state)

We've lost the FD number street city state \rightarrow zip

If we put this table back in the design, we are back where we started. And we violate BCNF.

65

© Lois Delcambre, David Maier 2005-2018

Algorithm for lossless join decomposition into BCNF relations (not necessarily dependency preserving)

••••
•••
•••
• •

66

set D := { R } (the current set of relations)
 while there is a relation in R that is not in BCNF relative

to FDs F

begin choose a relation Q that is not in BCNF find FD $X \rightarrow Y$ from F⁺ in Q that violates BCNF replace Q in D by two relations: (Q - Y) and (X U Y)

Finding a "troublesome" FD wot migne