

## Introduction to Database Systems CSE 444

### Lectures 8 & 9 Database Design

October 12 & 15, 2007

1

## Announcements/Reminders

- Homework 1: solutions are posted
- Homework 2: posted later today (due Sat., Oct. 20)
- Project Phase 1 due tomorrow, 9pm

2

## Outline

- The relational data model: 3.1
- Functional dependencies: 3.4

3

## Schema Refinements = Normal Forms

- 1st Normal Form = all tables are flat
- 2nd Normal Form = obsolete
- Boyce Codd Normal Form = will study
- 3rd Normal Form = see book

4

## First Normal Form (1NF)

- A database schema is in First Normal Form if all tables are flat

**Student**

Name	GPA	Courses
Alice	3.8	Math DB OS
Bob	3.7	DB OS
Carol	3.9	Math OS

→

May need to add keys

**Student**

Name	GPA
Alice	3.8
Bob	3.7
Carol	3.9

**Takes**

Student	Course
Alice	Math
Carol	Math
Alice	DB
Bob	DB
Alice	OS
Carol	OS

**Course**

Course
Math
DB
OS

5

## Relational Schema Design

Conceptual Model:

---

Relational Model:  
plus FD's

---

Normalization:  
Eliminates anomalies

6

## Data Anomalies

When a database is poorly designed we get anomalies:

**Redundancy:** data is repeated

**Update anomalies:** need to change in several places

**Delete anomalies:** may lose data when we don't want

7

## Relational Schema Design

Recall set attributes (persons with several phones):

Name	SSN	PhoneNumber	City
Fred	123-45-6789	206-555-1234	Seattle
Fred	123-45-6789	206-555-6543	Seattle
Joe	987-65-4321	908-555-2121	Westfield

One person may have multiple phones, but lives in only one city

**Anomalies:**

- Redundancy = repeated data
- Update anomalies = Fred moves to "Bellevue"
- Deletion anomalies = Joe deletes his phone number: what is his city ?

8

## Relation Decomposition

Break the relation into two:

Name	SSN	PhoneNumber	City
Fred	123-45-6789	206-555-1234	Seattle
Fred	123-45-6789	206-555-6543	Seattle
Joe	987-65-4321	908-555-2121	Westfield

Name	SSN	City	SSN	PhoneNumber
Fred	123-45-6789	Seattle	123-45-6789	206-555-1234
Joe	987-65-4321	Westfield	123-45-6789	206-555-6543
			987-65-4321	908-555-2121

**Anomalies are gone:**

- No more repeated data
- Easy to move Fred to "Bellevue" (how?)
- Easy to delete all Joe's phone numbers (how?)

9

## Relational Schema Design (or Logical Design)

Main idea:

- Start with some relational schema
- Find out its **functional dependencies**
- Use them to design a better relational schema

10

## Functional Dependencies

- A form of constraint
  - hence, part of the schema
- Finding them is part of the database design
- Also used in normalizing the relations

11

## Functional Dependencies

Definition:

If two tuples agree on the attributes

$$A_1, A_2, \dots, A_n$$

then they must also agree on the attributes

$$B_1, B_2, \dots, B_m$$

Formally:

$$A_1, A_2, \dots, A_n \rightarrow B_1, B_2, \dots, B_m$$

12

### When Does an FD Hold

Definition:  $A_1, \dots, A_m \rightarrow B_1, \dots, B_n$  holds in R if:

$$\forall t, t' \in R, (t.A_1=t'.A_1 \wedge \dots \wedge t.A_m=t'.A_m \Rightarrow t.B_1=t'.B_1 \wedge \dots \wedge t.B_n=t'.B_n)$$

R

	$A_1$	...	$A_m$		$B_1$	...	$B_n$	
t								
t'								

if t, t' agree here
then t, t' agree here

13

### Examples

An FD holds, or does not hold on an instance:

EmpID	Name	Phone	Position
E0045	Smith	1234	Clerk
E3542	Mike	9876	Salesrep
E1111	Smith	9876	Salesrep
E9999	Mary	1234	Lawyer

EmpID  $\rightarrow$  Name, Phone, Position  
 Position  $\rightarrow$  Phone  
 but not Phone  $\rightarrow$  Position

14

### Example

EmpID	Name	Phone	Position
E0045	Smith	1234	Clerk
E3542	Mike	9876	Salesrep
E1111	Smith	9876	Salesrep
E9999	Mary	1234	Lawyer

Position  $\rightarrow$  Phone

15

### Example

EmpID	Name	Phone	Position
E0045	Smith	1234	$\rightarrow$ Clerk
E3542	Mike	9876	Salesrep
E1111	Smith	9876	Salesrep
E9999	Mary	1234	$\rightarrow$ Lawyer

but not Phone  $\rightarrow$  Position

16

### Example

FD's are constraints:

- On some instances they hold
- On others they don't

$name \rightarrow color$   
 $category \rightarrow department$   
 $color, category \rightarrow price$

name	category	color	department	price
Gizmo	Gadget	Green	Toys	49
Tweaker	Gadget	Green	Toys	99

Does this instance satisfy all the FDs ?

17

### Example

$name \rightarrow color$   
 $category \rightarrow department$   
 $color, category \rightarrow price$

name	category	color	department	price
Gizmo	Gadget	Green	Toys	49
Tweaker	Gadget	Black	Toys	99
Gizmo	Stationary	Green	Office-sup.	59

What about this one ?

18

### An Interesting Observation

If all these FDs are true:

name	→	color
category	→	department
color, category	→	price

Then this FD also holds:

name, category	→	price
----------------	---	-------

Why ??

19

### Goal: Find ALL Functional Dependencies

- Anomalies occur when certain “bad” FDs hold
- We know some of the FDs
- Need to find *all* FDs, then look for the bad ones

20

### Armstrong’s Rules (1/3)

$A_1, A_2, \dots, A_n \rightarrow B_1, B_2, \dots, B_m$

Is equivalent to

$A_1, A_2, \dots, A_n \rightarrow B_1$   
 $A_1, A_2, \dots, A_n \rightarrow B_2$   
 . . . . .  
 $A_1, A_2, \dots, A_n \rightarrow B_m$

**Splitting rule  
and  
Combing rule**

	A1		Am		B1		Bm	

21

### Armstrong’s Rules (1/3)

$A_1, A_2, \dots, A_n \rightarrow A_i$

**Trivial Rule**

where  $i = 1, 2, \dots, n$

Why ?

	A1	...	Am	

22

### Armstrong’s Rules (1/3)

**Transitive Closure Rule**

If  $A_1, A_2, \dots, A_n \rightarrow B_1, B_2, \dots, B_m$

and  $B_1, B_2, \dots, B_m \rightarrow C_1, C_2, \dots, C_p$

then  $A_1, A_2, \dots, A_n \rightarrow C_1, C_2, \dots, C_p$

Why ?

23

	A1	...	Am		B1	...	Bm		C1	...	Cp	

24

### Example (continued)

Start from the following FDs:

1. name → color  
 2. category → department  
 3. color, category → price

Infer the following FDs:

Inferred FD	Which Rule did we apply ?
4. name, category → name	
5. name, category → color	
6. name, category → category	
7. name, category → color, category	
8. name, category → price	

25

### Example (continued)

Answers:

1. name → color  
 2. category → department  
 3. color, category → price

Inferred FD	Which Rule did we apply ?
4. name, category → name	Trivial rule
5. name, category → color	Transitivity on 4, 1
6. name, category → category	Trivial rule
7. name, category → color, category	Split/combine on 5, 6
8. name, category → price	Transitivity on 3, 7

THIS IS TOO HARD ! Let's see an easier way.

26

### Closure of a set of Attributes

**Given** a set of attributes  $A_1, \dots, A_n$

The **closure**,  $\{A_1, \dots, A_n\}^+$  = the set of attributes B s.t.  $A_1, \dots, A_n \rightarrow B$

Example:

name → color  
 category → department  
 color, category → price

Closures:

name<sup>+</sup> = {name, color}  
 {name, category}<sup>+</sup> = {name, category, color, department, price}  
 color<sup>+</sup> = {color}

27

### Closure Algorithm

$X = \{A_1, \dots, A_n\}$ .

**Repeat until X doesn't change do:**

if  $B_1, \dots, B_n \rightarrow C$  is a FD and  $B_1, \dots, B_n$  are all in X then add C to X.

Example:

name → color  
 category → department  
 color, category → price

$\{name, category\}^+ =$   
 { name, category, color, department, price }

Hence: name, category → color, department, price

28

### Example

In class:

R(A,B,C,D,E,F)

A, B → C  
 A, D → E  
 B → D  
 A, F → B

Compute  $\{A,B\}^+$  X = {A, B, }

Compute  $\{A, F\}^+$  X = {A, F, }

29

### Why Do We Need Closure

- With closure we can find all FD's easily
- To check if  $X \rightarrow A$ 
  - Compute  $X^+$
  - Check if  $A \in X^+$

30

## Using Closure to Infer ALL FDs

Example:

A, B	→	C
A, D	→	B
B	→	D

Step 1: Compute  $X^+$ , for every  $X$ :

$A^+ = A, B^+ = BD, C^+ = C, D^+ = D$
$AB^+ = ABCD, AC^+ = AC, AD^+ = ABCD,$
$BC^+ = BCD, BD^+ = BD, CD^+ = CD$
$ABC^+ = ABD^+ = ACD^+ = ABCD$ (no need to compute-- why ?)
$BCD^+ = BCD, ABCD^+ = ABCD$

Step 2: Enumerate all FD's  $X \rightarrow Y$ , s.t.  $Y \subseteq X^+$  and  $X \cap Y = \emptyset$ :

$AB \rightarrow CD, AD \rightarrow BC, ABC \rightarrow D, ABD \rightarrow C, ACD \rightarrow B$
---

31

## Another Example

- Enrollment(student, major, course, room, time)
  - student → major
  - major, course → room
  - course → time

What else can we infer ? [in class, or at home]

32

## Keys

- A **superkey** is a set of attributes  $A_1, \dots, A_n$  s.t. for any other attribute  $B$ , we have  $A_1, \dots, A_n \rightarrow B$
- A **key** is a minimal superkey
  - I.e. set of attributes which is a superkey and for which no subset is a superkey

33

## Computing (Super)Keys

- Compute  $X^+$  for all sets  $X$
- If  $X^+ =$  all attributes, then  $X$  is a key
- List only the minimal  $X$ 's

34

## Example

Product(name, price, category, color)

name, category	→	price
category	→	color

What is the key ?

35

## Example

Product(name, price, category, color)

name, category	→	price
category	→	color

What is the key ?

$(name, category)^+ = name, category, price, color$

Hence  $(name, category)$  is a key

36

### Examples of Keys

Enrollment(student, address, course, room, time)

student → address
room, time → course
student, course → room, time

(find keys at home)

37

### Eliminating Anomalies

Main idea:

- $X \rightarrow A$  is OK if X is a (super)key
- $X \rightarrow A$  is not OK otherwise

38

### Example

Name	SSN	PhoneNumber	City
Fred	123-45-6789	206-555-1234	Seattle
Fred	123-45-6789	206-555-6543	Seattle
Joe	987-65-4321	908-555-2121	Westfield
Joe	987-65-4321	908-555-1234	Westfield

$SSN \rightarrow Name, City$

What the key?

{SSN, PhoneNumber}

Hence  $SSN \rightarrow Name, City$  is a "bad" dependency

39

### Key or Keys ?

Can we have more than one key ?

Given  $R(A,B,C)$  define FD's s.t. there are two or more keys

40

### Key or Keys ?

Can we have more than one key ?

Given  $R(A,B,C)$  define FD's s.t. there are two or more keys

$AB \rightarrow C$	or	$A \rightarrow BC$
$BC \rightarrow A$		$B \rightarrow AC$

what are the keys here ?

Can you design FDs such that there are *three* keys ?

41

### Boyce-Codd Normal Form

A simple condition for removing anomalies from relations:

A relation R is in BCNF if:

If $A_1, \dots, A_n \rightarrow B$ is a non-trivial dependency in R, then $\{A_1, \dots, A_n\}$ is a superkey for R
---

In other words: there are no "bad" FDs

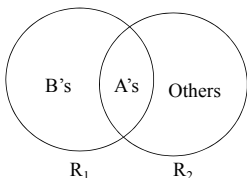
Equivalently:

$\forall X$ , either  $(X^+ = X)$  or  $(X^+ = \text{all attributes})$

42

## BCNF Decomposition Algorithm

**repeat**  
 choose  $A_1, \dots, A_m \rightarrow B_1, \dots, B_n$  that violates BCNF  
 split R into  $R_1(A_1, \dots, A_m, B_1, \dots, B_n)$  and  $R_2(A_1, \dots, A_m, [\text{others}])$   
 continue with both  $R_1$  and  $R_2$   
**until** no more violations



Is there a 2-attribute relation that is not in BCNF?

In practice, we have a better algorithm (coming up) <sup>43</sup>

## Example

Name	SSN	PhoneNumber	City
Fred	123-45-6789	206-555-1234	Seattle
Fred	123-45-6789	206-555-6543	Seattle
Joe	987-65-4321	908-555-2121	Westfield
Joe	987-65-4321	908-555-1234	Westfield

SSN  $\rightarrow$  Name, City

What the key?  
 {SSN, PhoneNumber} use SSN  $\rightarrow$  Name, City to split <sup>44</sup>

## Example

Name	SSN	City
Fred	123-45-6789	Seattle
Joe	987-65-4321	Westfield

SSN  $\rightarrow$  Name, City

SSN	PhoneNumber
123-45-6789	206-555-1234
123-45-6789	206-555-6543
987-65-4321	908-555-2121
987-65-4321	908-555-1234

Let's check anomalies:  
 • Redundancy?  
 • Update?  
 • Delete?

<sup>45</sup>

## Example Decomposition

Person(name, SSN, age, hairColor, phoneNumber)  
 SSN  $\rightarrow$  name, age  
 age  $\rightarrow$  hairColor

Decompose in BCNF (in class):

<sup>46</sup>

## BCNF Decomposition Algorithm

```

BCNF_Decompose(R)

  find X s.t.: X  $\neq X^+ \neq$  [all attributes]

  if (not found) then "R is in BCNF"

  let Y = X+ - X
  let Z = [all attributes] - X+
  decompose R into R1(X  $\cup$  Y) and R2(X  $\cup$  Z)
  continue to decompose recursively R1 and R2
    
```

<sup>47</sup>

Find X s.t.: X  $\neq X^+ \neq$  [all attributes]

## Example BCNF Decomposition

Person(name, SSN, age, hairColor, phoneNumber)  
 SSN  $\rightarrow$  name, age  
 age  $\rightarrow$  hairColor

Iteration 1: Person  
 SSN<sup>+</sup> = SSN, name, age, hairColor  
 Decompose into: P(SSN, name, age, hairColor)  
 Phone(SSN, phoneNumber)

Iteration 2: P  
 age<sup>+</sup> = age, hairColor  
 Decompose: People(SSN, name, age)  
 Hair(age, hairColor)  
 Phone(SSN, phoneNumber)

What are the keys?

<sup>48</sup>



R(A,B,C,D)

A → B  
B → C

### Example

R(A,B,C,D)  
A<sup>+</sup> = ABC ≠ ABCD

R<sub>1</sub>(A,B,C)  
B<sup>+</sup> = BC ≠ ABC

R<sub>2</sub>(A,D)

R<sub>11</sub>(B,C)

R<sub>12</sub>(A,B)

What are the keys ?

What happens if in R we first pick B<sup>+</sup> ? Or AB<sup>+</sup> ?

### Decompositions in General

R(A<sub>1</sub>, ..., A<sub>n</sub>, B<sub>1</sub>, ..., B<sub>m</sub>, C<sub>1</sub>, ..., C<sub>p</sub>)

R<sub>1</sub>(A<sub>1</sub>, ..., A<sub>n</sub>, B<sub>1</sub>, ..., B<sub>m</sub>)

R<sub>2</sub>(A<sub>1</sub>, ..., A<sub>n</sub>, C<sub>1</sub>, ..., C<sub>p</sub>)

R<sub>1</sub> = projection of R on A<sub>1</sub>, ..., A<sub>n</sub>, B<sub>1</sub>, ..., B<sub>m</sub>  
 R<sub>2</sub> = projection of R on A<sub>1</sub>, ..., A<sub>n</sub>, C<sub>1</sub>, ..., C<sub>p</sub>

### Theory of Decomposition

- Sometimes it is correct:

Name	Price	Category
Gizmo	19.99	Gadget
OneClick	24.99	Camera
Gizmo	19.99	Camera

Name	Price
Gizmo	19.99
OneClick	24.99
Gizmo	19.99

Name	Category
Gizmo	Gadget
OneClick	Camera
Gizmo	Camera

Lossless decomposition 51

### Incorrect Decomposition

- Sometimes it is not:

Name	Price	Category
Gizmo	19.99	Gadget
OneClick	24.99	Camera
Gizmo	19.99	Camera

Name	Category
Gizmo	Gadget
OneClick	Camera
Gizmo	Camera

Price	Category
19.99	Gadget
24.99	Camera
19.99	Camera

What's incorrect ??

Lossy decomposition 52

### Decompositions in General

R(A<sub>1</sub>, ..., A<sub>n</sub>, B<sub>1</sub>, ..., B<sub>m</sub>, C<sub>1</sub>, ..., C<sub>p</sub>)

R<sub>1</sub>(A<sub>1</sub>, ..., A<sub>n</sub>, B<sub>1</sub>, ..., B<sub>m</sub>)

R<sub>2</sub>(A<sub>1</sub>, ..., A<sub>n</sub>, C<sub>1</sub>, ..., C<sub>p</sub>)

If A<sub>1</sub>, ..., A<sub>n</sub> → B<sub>1</sub>, ..., B<sub>m</sub>  
Then the decomposition is lossless

Note: don't need A<sub>1</sub>, ..., A<sub>n</sub> → C<sub>1</sub>, ..., C<sub>p</sub>

BCNF decomposition is always lossless. WHY ? 53