

Relational Database Design(3)

Relational Database Design

- Introduction
- First Normal Form
- Pitfalls in Relational Database Design
- Functional Dependencies
- Decomposition
- Boyce-Codd Normal Form
- **Third Normal Form**
- **Multivalued Dependencies and Fourth Normal Form**
- **Overall Database Design Process**

Third Normal Form: Motivation

- There are some situations where
 - ❖ **BCNF** is not dependency preserving, and
 - ❖ **efficient checking for FD violation on updates is important**

Third Normal Form: Motivation

Solution: define a weaker normal form, called 3NF.

- ❖ Allows some redundancy (with resultant problems; we will see examples later)
- ❖ But FDs can be checked on individual relations without computing a join.
- ❖ ***There is always a lossless-join, dependency-preserving decomposition into 3NF.***

Algorithm

Third Normal Form (revisited)

- A relation schema R is in 3NF if for all:

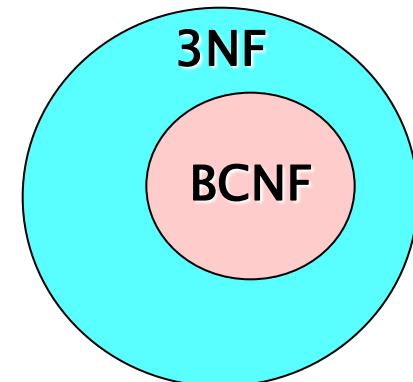
$$X \rightarrow Y \text{ in } F^+$$

at least one of the following holds:

- BCNF {
1. $X \rightarrow Y$ is trivial (i.e., $Y \in X$) or
 2. **X is a superkey for R or**
 3. If X is not a super key then,

for each attribute A in Y , $A \subseteq$ a candidate key for R .

Key $\rightarrow X \rightarrow K_1$ Reflexive, not transitive FD



(NOTE: each attribute may be in a different candidate key)

- Third condition is a minimal relaxation of BCNF to ensure dependency preservation (will see why later).

Third Normal Form (revisited)

$X \rightarrow Y$ in F^+

If X is not a super key then,

for each attribute A in Y , $A \subseteq$ a candidate key for R .

Key $\rightarrow X \rightarrow K_1$: X is not a key, but any key $\rightarrow X$

Key $\rightarrow K_1$: k_1 in Key so it is reflexive, not transitive

it is trivial, always true

Ex. Fname+Lname \rightarrow Lname

(NOTE: each attribute may be in a different candidate key)

Enrollment

Registration number: 205

Student : Mr. John Smith ID: 41010703
Department : Computer Science

CourseNo	Course Name	Unit	Section
729101	Accounting	4	700
729111	Finance	3	711
999211	Database	3	713

Total Units: 10
Total amount :5000

Not BCNF but in 3NF

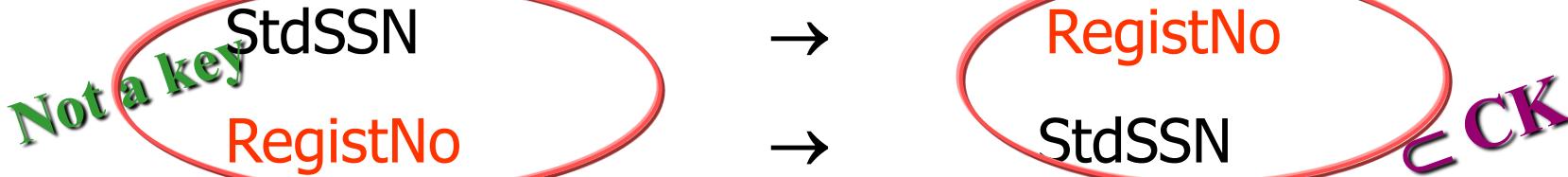
StdSSN	RegistNo	CourseNo	Section
41010703	205	729101	700
41010703	205	729111	711
41010703	205	999211	713
41010943	1368	729111	712
41010943	1368	999211	711
41010943	1368	729104	700
41012147	1684	729111	711
41012147	1684	999211	713

CK (StdSSN, CourseNo)

→ RegistNo, Section

CK (RegistNo, CourseNo)

→ StdSSN, Section



CK (StdSSN, CourseNo) → RegistNo, Section

RegistNo
Not a key

→ StdSSN
CK

(StdSSN, CourseNo) → RegistFormNo → StdSSN

Hence (StdSSN, CourseNo) → StdSSN

reflexive, always true

3NF

**allows a transitive relationship
such that**

key → y → k₁

a part of key

Canonical Cover

- Sets of FD may have redundant dependencies that can be inferred from the others

❖ Eg: $A \rightarrow C$ is redundant in:

$\{A \rightarrow B, B \rightarrow C, A \rightarrow C\}$

$A \rightarrow C$ can be inferred from $\{A \rightarrow B, B \rightarrow C\}$

$\{A \rightarrow B, B \rightarrow C, A \rightarrow C\}$ equivalent to

$\{A \rightarrow B, B \rightarrow C\}$

Canonical Cover

- Parts of a FD may be redundant

On RHS: $\{A \rightarrow B, B \rightarrow C, A \rightarrow CD\}$

can be simplified to

$\{A \rightarrow B, B \rightarrow C, A \rightarrow D\}$

because

$A \rightarrow CD$ means $A \rightarrow C$ and $A \rightarrow D$

Canonical Cover

- Parts of a FD may be redundant

On LHS: $\{A \rightarrow B, B \rightarrow C, AC \rightarrow D\}$

can be simplified to

$\{A \rightarrow B, B \rightarrow C, A \rightarrow D\}$

because

$A \rightarrow C$ hence we can drop C on LHS of $AC \rightarrow D$

Ex. $(ID, Name) \rightarrow Telephone$ equi to

$ID \rightarrow Telephone$

because $ID \rightarrow Name$

Canonical Cover

- canonical cover of F (F_c) is a “minimal” set of FDs equivalent to F , having no redundant dependencies or redundant parts of dependencies

F_c is equivalent to F

3NF Decomposition Algorithm

Find a canonical cover (F_c) for F ;

$R_c = \{ \}$

for each functional dependency $X \rightarrow Y$ in F_c do

{ create a schema XY if $XY \not\subset$ any R_i

$R_c = R_c \cup \{XY\}$

}

if none of the schemas in R_c

contains a candidate key for R

then $R_c = R_c \cup \{CK\}$

where CK is any candidate key for R

return R_c

3NF Decomposition Algorithm (Cont.)

■ Above algorithm ensures:

- ❖ each relation schema R_i is in 3NF**
- ❖ decomposition is dependency preserving and lossless-join**

Example

- Relation schema:

(A, B, C, D)

- The functional dependencies for this relation schema are:

$C \rightarrow A, D$

$A, B \rightarrow C$

Not in 3NF

Has transitive FD

- The key is: $\{A, B\}$

Applying 3NF to Banker-info-schema

$$F_C = \{ C \rightarrow (A, D), (A, B) \rightarrow C \}$$

- In the **for** loop, we add :

$$R1 = (A, \underline{C}, D)$$

$$R2 = (\underline{A}, B, C)$$

- Since $R2$ contains a candidate key for $R1$, we are done with the decomposition process.
- Return $R_C = \{R1, R2\}$

3NF Decomposition Algorithm

1. Find **keys**, and check if the relation is already in **3NF**

2. **Decompose** the right hand sides of the dependencies

$$X \rightarrow YZ \implies X \rightarrow Y, X \rightarrow Z$$

3. Remove redundant attributes on the left hand sides

$$XY \rightarrow Y \implies X \rightarrow Y$$

4. Remove redundant functional dependencies

5. **Combine** dependencies with same left hand sides

6. **Create a relation for each functional dependency**

7. Remove relations contained in other relations

8. If no relation contains a key of the original relation, add a relation whose attributes form such a key

1.Finding the Keys

$R = \{ A, B, C, D, E, F, G \}$

$(A \rightarrow BC), (B \rightarrow C), (AB \rightarrow D),$
 $(EF \rightarrow G), (G \rightarrow F)$

L = attributes appearing **only** on left sides = AE

R = attributes appearing **only** on right sides = CD

N = attributes **not** appearing at all = { }

Every key must include L and N

Every key must be disjoint from R

Hence, $\{A, E\} \subseteq \text{key} \subseteq \{ A, E, B, F, G \}$

1.Finding the Keys (cont'd)

$\{ A,E \} \subseteq \text{key} \subseteq \{ A,E,B,F,G \}$ appearing on left sides

$(A \rightarrow BC), (B \rightarrow C), (AB \rightarrow D),$

$(EF \rightarrow G), (G \rightarrow F)$

$\{ A,E \}^+ = A,E,B,C,D \neq R$

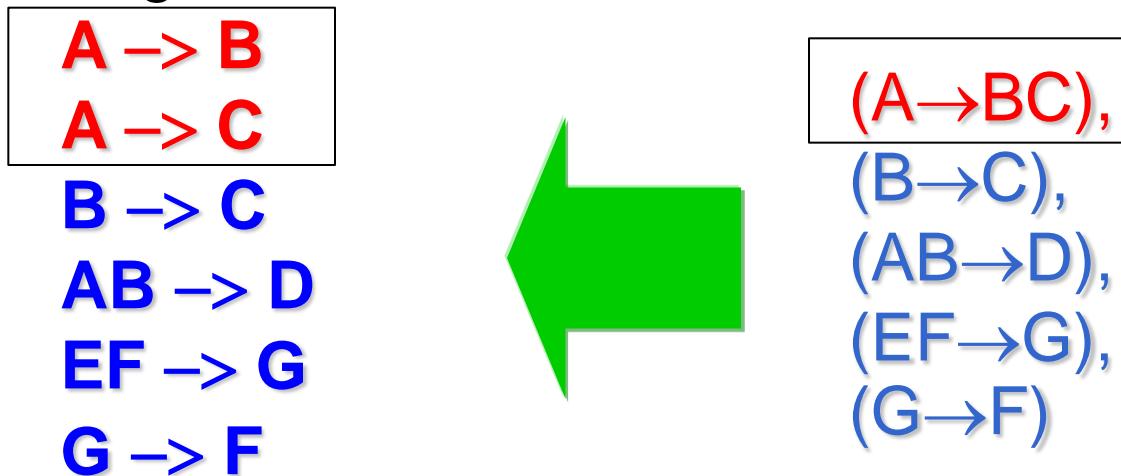
$\{ A,E,F \}^+ = A,E,B,C,D,F,G = R$

$\{ A,E,G \}^+ = A,E,B,C,D,G,F = R$

Hence, keys are $\{A,E,F\}$ and $\{A,E,G\}$

2. Decomposing right hand sides

- Decompose the right hand sides of the functional dependencies, e.g.,
 - $X \rightarrow YZ$ becomes $X \rightarrow Y$ and $X \rightarrow Z$.
- This gives us:



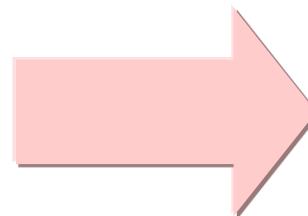
3. Remove redundant attributes

$(A \rightarrow BC)$, $(B \rightarrow C)$, $(AB \rightarrow D)$, $(EF \rightarrow G)$, $(G \rightarrow F)$

- Remove redundant attributes from the left sides
- Check $(AB \rightarrow D)$ for redundancy on left:
 - Since $B^+ = C$; A is needed to obtain D*
 - Since $A^+ = ABCD$; A suffices for D*
 - Hence, we can replace $AB \rightarrow D$ with $(A \rightarrow D)$*
- Check $(EF \rightarrow D)$ for redundancy on left:
 - Both E and F are needed to get D, so the rule cannot be replaced

The Simplified Rules

- 1. $A \rightarrow B$
- 2. $A \rightarrow C$
- 3. $B \rightarrow C$
- 4. $AB \rightarrow D$
- 5. $EF \rightarrow G$
- 6. $G \rightarrow F$



- 1. $A \rightarrow B$
- 2. $A \rightarrow C$
- 3. $B \rightarrow C$
- 4. $A \rightarrow D$
- 5. $EF \rightarrow G$
- 6. $G \rightarrow F$

4. Remove the Redundant Rules

- $(A \rightarrow B), (A \rightarrow C), (B \rightarrow C), (A \rightarrow D),$
 $(EF \rightarrow G), (G \rightarrow F)$
- Is $(A \rightarrow B)$ redundant?
 - ❖ Compute A^+ without using this rule :
 - ☞ $A^+ = ACD$
 - ☞ Since B not in ACD , this rule is **not redundant**
- Is $(A \rightarrow C)$ redundant
 - ❖ Compute A^+ without using this rule:
 - ☞ $A^+ = ABCD$
 - ☞ Redundant, since C is in A^+
- Remove the Redundant Rules $(A \rightarrow C)$.

4.Remove the Redundant Rules(Cont.)

(A → B), ~~(A → C)~~, (B → C), (A → D), (EF → G), (G → F)

- Is (B → C) redundant? Compute B+ :

- ❖ $B^+ = B$
 - ❖ Hence, not redundant

- Similarly, rules (A → D), (EF → G), and (G → F) are not redundant

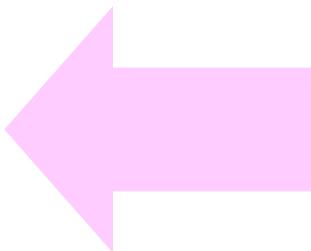
- ❖ Resulting rules, after renumbering:

1. A → B
2. B → C
3. A → D
4. EF → G
5. G → F

5. Combining Functional Dependencies

- Combine all dependencies with the same left side
- This gives:

1. $A \rightarrow BD$
2. $B \rightarrow C$
3. $EF \rightarrow G$
4. $G \rightarrow F$



1. $A \rightarrow B$
2. $B \rightarrow C$
3. $A \rightarrow D$
4. $EF \rightarrow G$
5. $G \rightarrow F$

6.Creating Relations

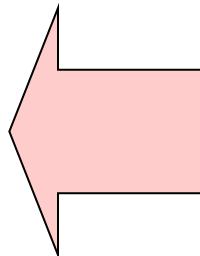
- Create a relation for each dependency
- This gives:

1. ABD

2. BC

3. EFG

4. GF



1. $A \rightarrow BD$

2. $B \rightarrow C$

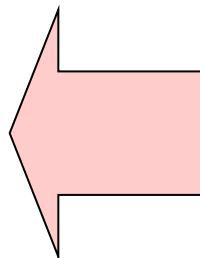
3. $EF \rightarrow G$

4. $G \rightarrow F$

7.Remove Redundant Relations

- Remove relations contained in other relations
- This gives:

1. ABD
2. BC
3. EFG



1. ABD
2. BC
3. EFG
4. GF

subset

8. Assuring Storage of a Global Key

- If no relation contains a key of the original relation, then add a relation whose attributes form such a key
- None of the relations contains **AEF or AEG (keys)**
 - ❖ Hence, one of these has to be added
- Two choices for the final decomposition into 3NF:
 1. **ABD, BC, EFG, AEF**
 2. **ABD, BC, EFG, AEG**

Why was the Global Key Needed?

- Consider $R(A,B,X,Y)$, with $(A \rightarrow X)$ and $(B \rightarrow Y)$
- The global key is AB
- Hence, the decomposition AX, BY is not lossless join
 - ❖ $AX \cap BY = \{ \}$
- But AX, BY, AB is in 3NF and lossless join+ dependency preserving
 - ❖ $AX \cap AB = A : A \rightarrow AX,$
 - ❖ $AB \cap BY = B : B \rightarrow BY$

AB, the key, joins all tables together

Comparison of BCNF and 3NF

- It is always possible to decompose a relation into relations in 3NF and
 - ❖ the decomposition is lossless
 - ❖ the dependencies are preserved
- It is always possible to decompose a relation into relations in BCNF and
 - ❖ the decomposition is lossless
 - ❖ *it may not be possible to preserve dependencies.*

3NF (Cont.)

■ Example

- ❖ $R = (A, B, C)$
 $F = \{AB \rightarrow C, C \rightarrow B\}$
- ❖ Two candidate keys: AB and AC
- ❖ R is in 3NF

$AB \rightarrow C$ AB is a superkey
 $C \rightarrow B$ B is contained in a candidate key

- BCNF decomposition has (AC) and (CB)
 - Testing for $AB \rightarrow C$ requires a join

- There is some redundancy in this schema
- Equivalent to example in book:

Banker-schema = (branch-name, customer-name, banker-name)
banker-name \rightarrow branch name
branch name ,customer-name \rightarrow banker-name

Design Goals

- Goal for a relational database design is:
 - ❖ BCNF.
 - ❖ Lossless join.
 - ❖ Dependency preservation.
- If we cannot achieve this, we accept one of
 - ❖ Lack of dependency preservation
 - ❖ Redundancy due to use of 3NF
- Interestingly, **SQL does not provide a direct way of specifying functional dependencies other than superkeys.**
Can specify FDs using assertions, but they are expensive to test
- Even if we had a dependency preserving decomposition, using SQL we would **not be able to efficiently test a functional dependency whose left hand side is not a key.**

Multivalued Dependencies

- Consider a database

classes(course, teacher, book)

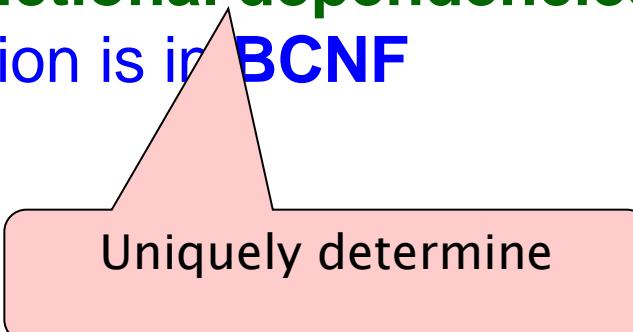
- *Teachers and books assigned in each course*

1 course - many teachers (course->>teacher)

1 course - many books (course->>book)

No relation between teachers and books

- no non-trivial functional dependencies and therefore the relation is in BCNF



Uniquely determine

Multivalued Dependencies (Cont.)

classes	course	teacher	book
	database	Avi	DB Concepts
	database	Avi	Ullman
	database	Hank	DB Concepts
	database	Hank	Ullman
	database	Sudarshan	DB Concepts
	database	Sudarshan	Ullman
	operating systems	Avi	OS Concepts
	operating systems	Avi	Shaw
	operating systems	Jim	OS Concepts
	operating systems	Jim	Shaw

repeat

- BCNF but there are redundant data.
- Insertion anomalies – i.e., if Sara is a new teacher that can teach database, two tuples need to be inserted
 - (database, Sara, DB Concepts)
 - (database, Sara, Ullman)

Multivalued Dependencies (Cont.)

- Therefore, it is better to decompose *classes* into:

course	teacher	
database	Avi	
database	Hank	
database	Sudarshan	
operating systems	Avi	
operating systems	Jim	

teaches

course	book	
database	DB Concepts	
database	Ullman	
operating systems	OS Concepts	
operating systems	Shaw	

text

We shall see that these two relations are in **Fourth Normal Form (4NF)**

Multivalued Dependencies

- The concept of multivalued dependencies was developed to provide a basis for decomposition of relations like the classes scheme.
- Definition

- $X \rightarrow\!> Y|Z$ is said to hold for $R(X, Y, Z)$ if t_1 and t_2 are given by

$t_1 = [X, Y_1, Z_1]$, and

$t_2 = [X, Y_2, Z_2]$

- then there must be tuples t_3 and t_4 such that

$t_3 = [X, Y_1, Z_2]$, and

$t_4 = [X, Y_2, Z_1]$

independent

course	teacher	book
database database	Avi Sara	DB Concepts Ullman
database database	Avi Sara	Ullman DB Concepts

MVD

- Dependency between attributes (for example, A, B, and C) in a relation, such that for each value of A there is a set of values for B and a set of values for C. However, set of values for B and C are independent of each other.

Fourth Normal Form (4NF)

- MVD between attributes A, B, and C in a relation using the following notation:

A $\rightarrow\!\!>$ B

A $\rightarrow\!\!>$ C

Fourth Normal Form (4NF)

- **Defined as a relation that is in BCNF and contains no nontrivial MVDs.**

No MVDs

Fourth Normal Form

- A relation schema R is in 4NF with respect to a set D if for all MVD in D^+ of the form $X \twoheadrightarrow Y$,
where $X \subseteq R$ and $Y \subseteq R$, **at least one of the following hold:**
 - ❖ $X \twoheadrightarrow Y$ is trivial (i.e., $X \subseteq Y$ or $\underline{Y \cup X = R}$) or
 - ❖ X is a superkey for schema R
- If a relation is in 4NF it is also in BCNF.
for a relation with at least three attributes having multi-value dependency

no MVDs or $\underline{Y \cup X = R}$

4NF Decomposition Algorithm

```
result := {R};  
done := false;  
compute  $D^+$ ;  
Let  $D_i$  denote the restriction of  $D^+$  to  $R_i$ 
```

Similar to BCNF
algorithm

while (not done)

if (there is a schema R_i in result that is not in 4NF) **then**
begin

let $Y \rightarrow\!\!\rightarrow X$ be a nontrivial multivalued dependency that holds
on R_i such that $Y \rightarrow\!\!\rightarrow R_i$ is not in D_i , and $Y \cap X = \emptyset$;
result := (result - R_i) \cup ($R_i - X$) \cup (Y, X); {Split into 2 schemas}

end

else done := true;

Note: each R_i is in 4NF, and decomposition is lossless-join



Example

- $R = (A, B, C, G, H, I)$

$$F = \{ A \rightarrow\!\!\!\rightarrow B \\ B \rightarrow\!\!\!\rightarrow HI \\ CG \rightarrow\!\!\!\rightarrow H \}$$

- R is not in 4NF since $A \rightarrow\!\!\!\rightarrow B$ and A is not a super key for R

- Decomposition

- a) $R_1 = (A, B)$ (R_1 is in 4NF; $X \cup Y = R_1$)
- b) $R_2 = (\textcolor{red}{A}, C, G, H, I)$ (R_2 is not in 4NF; $CG \rightarrow\!\!\!\rightarrow H$)
- c) $R_3 = (C, G, H)$ (R_3 is in 4NF)
- d) $R_4 = (A, \textcolor{red}{C}, \textcolor{red}{G}, I)$ (R_4 is not in 4NF)
- Since $A \rightarrow\!\!\!\rightarrow B$ and $B \rightarrow\!\!\!\rightarrow HI$, hence $A \rightarrow\!\!\!\rightarrow HI$, $A \rightarrow\!\!\!\rightarrow I$
- e) $R_5 = (A, I)$ (R_5 is in 4NF)
- f) $R_6 = (\textcolor{red}{A}, C, G)$ (R_6 is in 4NF)

รีลัชನ์สหารสомн-การปรึกษา

TeacherID	TcourseNo	SID
Q1011	999211	41010943
Q1011	999211	41013780
Q1035	729101	41012451
Q1035	729104	41012451
	999211	41012451
	729111	41010703
	729111	41013327
	729103	41012147

No relationship between course teaching and students

TeacherID->> TcourseNo
TeacherID->> SID

$\downarrow X->>Y \text{ and } X \cup Y = R \downarrow$

รีลัชន์สหสาขาวารป-นักศึกษา

TeacherID	TcourseNo
Q1011	999211
Q1035	729101
Q1035	729104
Q1035	999211
Q1059	729111
Q1061	729103

รีลัชន์สหสาขาวารป-นักเรียนต้อง

TeacherID	SID
Q1011	41010943
Q1011	41013780
Q1035	41012451
Q1059	41010703
Q1059	41013327
Q1061	41012147

Fifth Normal Form

- The aim of is to have relations that cannot be decomposed further.
- A relation R satisfies **join dependency** (R_1, R_2, \dots, R_n) if and only if R is equal to the join of R_1, R_2, \dots, R_n where $R_i \subseteq R$.
- A relation R is in **5NF** (or project-join normal form, PJNF) if for all join dependencies, at least one of the following holds.
 - (R_1, R_2, \dots, R_n) is a trivial join-dependency (that is, one of R_i is R) **or**
 - Every R_i is a candidate key for R .

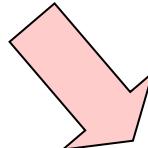
schema contains only key , can not be decomposed further

5NF

The diagram illustrates dependencies between the columns of the table. A red curved arrow points from the **dept** column to the **student** column. A green curved arrow points from the **subject** column to the **student** column. A red arrow points directly from the **subject** column to the **student** column.

<i>dept</i>	<i>subject</i>	<i>student</i>	not 5NF
<i>Comp. Sc.</i>	<i>CP1000</i>	<i>John Smith</i>	
<i>Mathematics</i>	<i>MA1000</i>	<i>John Smith</i>	
<i>Comp. Sc.</i>	<i>CP2000</i>	<i>Arun Kumar</i>	
<i>Comp. Sc.</i>	<i>CP3000</i>	<i>Reena Rani</i>	
<i>Physics</i>	<i>PH1000</i>	<i>Raymond Chew</i>	
<i>Chemistry</i>	<i>CH2000</i>	<i>Albert Garcia</i>	

no MVDs since subject and student are dependent;



5NF

(*dept*, *subject*)
 (*dept*, *student*)
 (*subject*, *student*)

5NF

Joining these 3 tables
gets the original
table

<u>Dept</u>	<u>Subject</u>
Comp. Sc.	CP1000
Mathematics	MA1000
Comp. Sc.	CP2000
Comp. Sc.	CP3000
Physics	PH1000
Chemistry	CH2000

<u>Dept</u>	<u>Student</u>
Comp. Sc.	John Smith
Mathematics	John Smith
Comp. Sc.	Arun Kumar
Comp. Sc.	Reena Rani
Physics	Raymond Chew
Chemistry	Albert Garcia

(dept, subject)
(dept, student)
(subject, student)

<u>Subject</u>	<u>Student</u>
CP1000	John Smith
MA1000	John Smith
CP2000	Arun Kumar
CP3000	Reena Rani
PH1000	Raymond Chew
CH2000	Albert Garcia

5NF

รากชั้นรายการ

รหัสสุจิสาขา	หมู่เรียน	ชื่ออาจารย์ผู้สอน
729101	700	ศิริชัย ศรีพราหมณ์
729111	711	จันทนา พราหมณ์ศิริ
729111	712	ศิริชัย ศรีพราหมณ์
729211	700	ศิริกฤติรา เหงื่อนมาดัย
729211	700	ศิริชัย ศรีพราหมณ์

not a key

รากชั้นที่ 1

รหัสสุจิสาขา	หมู่เรียน
729101	700
729111	711
729111	712
729211	700

รากชั้นที่ 2

หมู่เรียน	ชื่ออาจารย์ผู้สอน
700	ศิริชัย ศรีพราหมณ์
711	จันทนา พราหมณ์ศิริ
712	ศิริชัย ศรีพราหมณ์
700	ศิริกฤติรา เหงื่อนมาดัย

รากชั้นที่ 3

รหัสสุจิสาขา	ชื่ออาจารย์ผู้สอน
729101	ศิริชัย ศรีพราหมณ์
729111	จันทนา พราหมณ์ศิริ
729111	ศิริชัย ศรีพราหมณ์
729211	ศิริกฤติรา เหงื่อนมาดัย
729211	ศิริชัย ศรีพราหมณ์

Lossy Join

รากชั้นที่ 1 จัดจากการรวมรากชั้นที่ 1 รากชั้นที่ 2 และรากชั้นที่ 3

รหัสสุจิสาขา	หมู่เรียน	ชื่ออาจารย์ผู้สอน
729101	700	ศิริชัย ศรีพราหมณ์
729101	700	ศิริกฤติรา เหงื่อนมาดัย
729111	711	จันทนา พราหมณ์ศิริ
729111	712	ศิริชัย ศรีพราหมณ์
729211	700	ศิริชัย ศรีพราหมณ์
729211	700	ศิริกฤติรา เหงื่อนมาดัย

Spurious Tuple
more false data

Other Design Issues

■ Examples of bad database design, to be avoided:

Instead of *earnings(company-id, year, amount)*, use

- ❖ *earnings-2000, earnings-2001, earnings-2002, etc., all on the schema (company-id, earnings).*
 - Above are in BCNF, but make querying across years difficult and needs new table each year
- ❖ *company-year(company-id, earnings-2000, earnings-2001, earnings-2002)*
 - Also in BCNF, but also makes querying across years difficult and requires new attribute each year.
 - Is an example of a **crosstab**, where values for one attribute become column names
 - Used in spreadsheets, and in data analysis tools

Summarize

Ordinary Table $\xrightarrow{\text{minimizes}}$ 3NF , BCNF $\xrightarrow{\text{two}}$

Associated Table $\xrightarrow{\text{in}}$ 4NF , 5NF

Question ?

