# Decomposition & Schema Normalization

*CS 564- Fall 2016*

# HOW TO BUILD A DB APPLICATION

- Pick an application
- Figure out what to model (ER model)
  - Output: ER diagram
- Transform the ER diagram to a relational schema

- Refine the relational schema (normalization)

- Now ready to implement the schema and load the data!

# DB DESIGN THEORY

- Helps us identify the "bad" schemas and improve them
  1. express constraints on the data: functional dependencies (**FDs**)
  2. use the FDs to decompose the relations

- The process, called normalization, obtains a schema in a "normal form" that guarantees certain properties
  - examples of normal forms: **BCNF**, **3NF**, …

# SCHEMA DECOMPOSITION

# WHAT IS DECOMPOSITION?

We decompose $\mathbf{R}(A_1, ..., A_n)$ by creating

- $\mathbf{R_1}(B_1, .., B_m)$

- $\mathbf{R_2}(C_1, ..., C_l)$

- where $\{B_1, ..., B_m\} \cup \{C_1, ..., C_l\} = \{A_1, ... A_n\}$

- The instance of $\mathbf{R_1}$ is the projection of $\mathbf{R}$ onto $B_1, .., B_m$
- The instance of $\mathbf{R_2}$ is the projection of $\mathbf{R}$ onto $C_1, .., C_l$

# EXAMPLE: DECOMPOSITION

| SSN | name | age | phoneNumber |
|---|---|---|---|
| 934729837 | Paris | 24 | 608-374-8422 |
| 934729837 | Paris | 24 | 603-534-8399 |
| 123123645 | John | 30 | 608-321-1163 |
| 384475687 | Arun | 20 | 206-473-8221 |

| SSN | name | age |
|---|---|---|
| 934729837 | Paris | 24 |
| 123123645 | John | 30 |
| 384475687 | Arun | 20 |

| SSN | phoneNumber |
|---|---|
| 934729837 | 608-374-8422 |
| 934729837 | 603-534-8399 |
| 123123645 | 608-321-1163 |
| 384475687 | 206-473-8221 |

# DECOMPOSITION DESIDERATA

What should a good decomposition achieve?

1. minimize redundancy
2. avoid information loss (lossless-join)
3. preserve the FDs (dependency preserving)
4. ensure good query performance

# EXAMPLE: INFORMATION LOSS

| name | age | phoneNumber |
|------|-----|-------------|
| Paris | 24 | 608-374-8422 |
| John | 24 | 608-321-1163 |
| Arun | 20 | 206-473-8221 |

Decompose into:
**R₁**(name, age)
**R₂**(age, phoneNumber)

| name | age |
|------|-----|
| Paris | 24 |
| John | 24 |
| Arun | 20 |

| age | phoneNumber |
|-----|-------------|
| 24 | 608-374-8422 |
| 24 | 608-321-1163 |
| 20 | 206-473-8221 |

*Can we put it back together?*

# LOSSLESS-JOIN DECOMPOSITION

$\mathbf{R}$(A, B, C)

*decompose (projection)*

$\mathbf{R_1}$(A, B)     $\mathbf{R_2}$(B, C)

*recover (natural join)*

$\mathbf{R'}$(A, B, C)

A schema decomposition is **lossless-join** if for any initial instance $\mathbf{R}$, $\mathbf{R}$ = $\mathbf{R'}$

# LOSSLESS-JOIN CRITERION

- relation **R**(**A**) + set $F$ of FDs

- decomposition of **R** into $\mathbf{R_1(A_1)}$ and $\mathbf{R_2(A_2)}$

A decomposition is lossless-join **if and only if** at least one of the FDs is in $F^+$ (the closure of $F$) :

1. $A_1 \cap A_2 \longrightarrow A_1$
2. $A_1 \cap A_2 \longrightarrow A_2$

# EXAMPLE

- relation **R**(A, B, C, D)
- FD  $A \longrightarrow B, C$

lossless-join

- decomposition into **R$_1$**(A, B, C) and **R$_2$**(A, D)

**Not** lossless-join

- decomposition into **R$_1$**(A, B, C) and **R$_2$**(D)

# DEPENDENCY PRESERVING

Given **R** and a set of FDs $F$, we decompose **R** into **R$_1$** and **R$_2$**. Suppose:

- **R$_1$** has a set of FDs $F_1$
- **R$_2$** has a set of FDs $F_2$
- $F_1$ and $F_2$ are computed from $F$

A decomposition is **<u>dependency preserving</u>** if by enforcing $F_1$ over **R$_1$** and $F_2$ over **R$_2$**, we can enforce $F$ over **R**

# GOOD EXAMPLE

**Person**(SSN, name, age, canDrink)

- $SSN \longrightarrow name, age$

- $age \longrightarrow canDrink$

decomposes into

- **R$_1$**(SSN, name, age)
  - $SSN \longrightarrow name, age$
- **R$_2$**(age, canDrink)
  - $age \longrightarrow canDrink$

# BAD EXAMPLE

**R**(A, B, C)

- $A \longrightarrow B$

- $B, C \longrightarrow A$

Decomposes into:

- **R₁**(A, B)

  - $A \longrightarrow B$

- **R₂**(A, C)

  - no FDs here!!

**R₁**

| A | B |
|---|---|
| $a_1$ | b |
| $a_2$ | b |

**R₂**

| A | C |
|---|---|
| $a_1$ | c |
| $a_2$ | c |

*recover*

| A | B | C |
|---|---|---|
| $a_1$ | b | c |
| $a_2$ | b | c |

The recovered table violates $B, C \longrightarrow A$

# Normal Forms

A **<u>normal form</u>** represents a "good" schema design:

- 1NF (flat tables/atomic values)
- 2NF
- **3NF**
- **BCNF**
- 4NF
- ...

<span style="color:darkred">more restrictive</span>

# BCNF DECOMPOSITION

# BOYCE-CODD NORMAL FORM (BCNF)

A relation **R** is in **<u>BCNF</u>** if whenever $X \longrightarrow B$ is a non-trivial FD, then $X$ is a <span style="color:red">superkey</span> in **R**

**Equivalent definition**: for every attribute set $X$

- either $X^+ = X$
- or $X^+ = all\ attributes$

# BCNF EXAMPLE 1

| SSN | name | age | phoneNumber |
|---|---|---|---|
| 934729837 | Paris | 24 | 608-374-8422 |
| 934729837 | Paris | 24 | 603-534-8399 |
| 123123645 | John | 30 | 608-321-1163 |
| 384475687 | Arun | 20 | 206-473-8221 |

$$SSN \longrightarrow name, age$$

- **key** $= \{SSN, phoneNumber\}$
- $SSN \longrightarrow name, age$ is a "bad" FD
- The above relation is **not** in BCNF!

# BCNF Example 2

| SSN | name | age |
|-----|------|-----|
| 934729837 | Paris | 24 |
| 123123645 | John | 30 |
| 384475687 | Arun | 20 |

$$SSN \longrightarrow name, age$$

- **key** = $\{SSN\}$
- The above relation is in BCNF!

# BCNF EXAMPLE 3

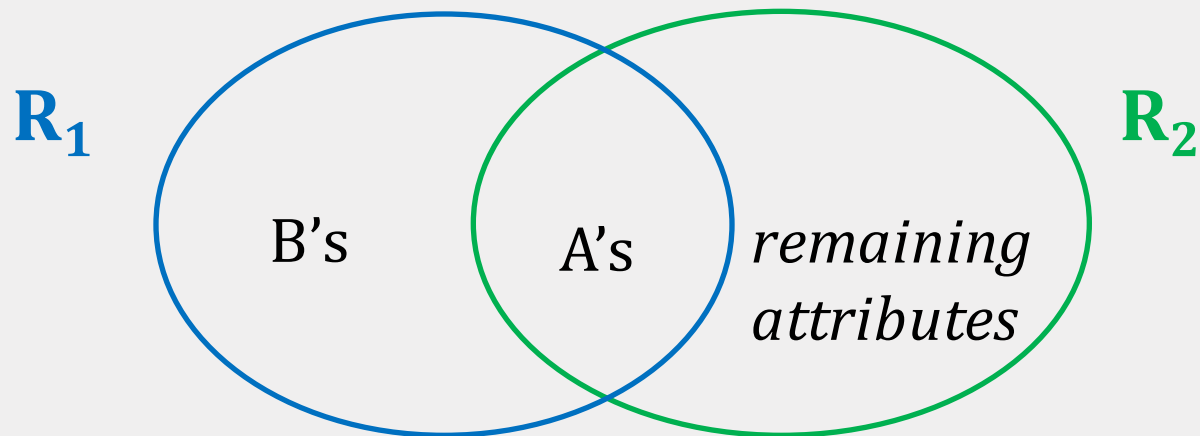| SSN | phoneNumber |
|-----|-------------|
| 934729837 | 608-374-8422 |
| 934729837 | 603-534-8399 |
| 123123645 | 608-321-1163 |
| 384475687 | 206-473-8221 |

- **key** $= \{SSN, phoneNumber\}$
- The above relation is in BCNF!
- **Q**: is it possible that a binary relation is not in BCNF?

# BCNF DECOMPOSITION

- Find an FD that violates the BCNF condition

$$A_1, A_2, ..., A_n \longrightarrow B_1, B_2, ..., B_m$$

- Decompose **R** to $\mathbf{R_1}$ and $\mathbf{R_2}$:
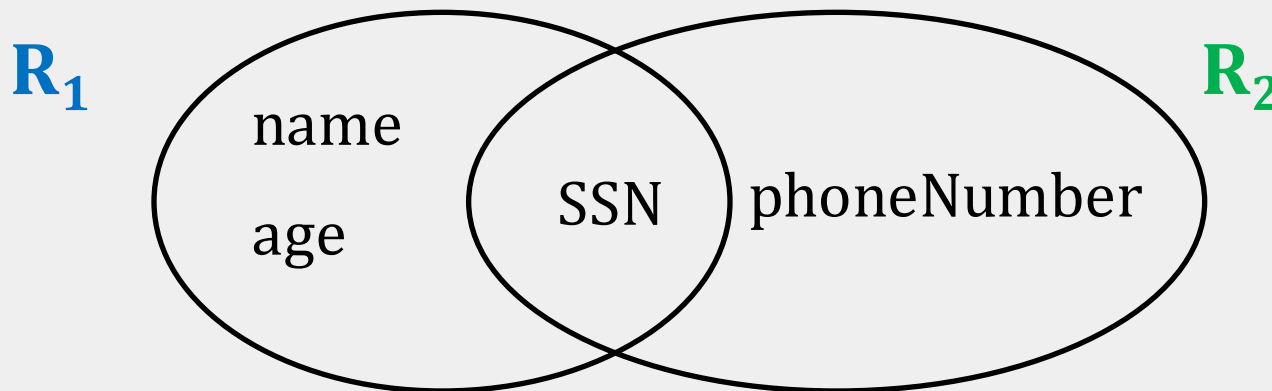
$\mathbf{R_1}$                  $\mathbf{R_2}$

B's     A's   *remaining attributes*

- Continue until no BCNF violations are left

# EXAMPLE

| SSN | name | age | phoneNumber |
|---|---|---|---|
| 934729837 | Paris | 24 | 608-374-8422 |
| 934729837 | Paris | 24 | 603-534-8399 |
| 123123645 | John | 30 | 608-321-1163 |
| 384475687 | Arun | 20 | 206-473-8221 |

- The FD $SSN \longrightarrow name, age$ violates BCNF

- Split into two relations $R_1$, $R_2$ as follows:

$R_1$ $R_2$

name
age
SSN
phoneNumber

# EXAMPLE CONT'D

**R₁**

name

age

SSN

phoneNumber

**R₂**

$$SSN \longrightarrow name, age$$

| SSN | name | age |
|-----|------|-----|
| 934729837 | Paris | 24 |
| 123123645 | John | 30 |
| 384475687 | Arun | 20 |

| SSN | phoneNumber |
|-----|-------------|
| 934729837 | 608-374-8422 |
| 934729837 | 603-534-8399 |
| 123123645 | 608-321-1163 |
| 384475687 | 206-473-8221 |

# BCNF DECOMPOSITION PROPERTIES

BCNF decomposition:

- – removes certain types of redundancy
- – is lossless-join
- – is not always dependency preserving

# BCNF is Lossless-Join

Example:

$R(A, B, C)$ with $A \longrightarrow B$ decomposes into:

$R_1(A, B)$ and $R_2(A, C)$

- BCNF decomposition satisfies the lossless-join criterion!

# BCNF is Not Dependency Preserving

**R**(A, B, C)

- $A \longrightarrow B$

- $B, C \longrightarrow A$

The BCNF decomposition is:

- **R$_1$**(A, B) with FD $A \longrightarrow B$
- **R$_2$**(A, C) with no FDs

There may not exist any BCNF decomposition that is FD preserving!

# BCNF Example (1)

**Books** (author, gender, booktitle, genre, price)

- *author → gender*

- *booktitle → genre, price*

What is the candidate key?

- *(author, booktitle)* is the only one!

Is is in BCNF?

- **No**, because the left hand side of both (not trivial) FDs is not a superkey!

# BCNF EXAMPLE (2)

**Books** (author, gender, booktitle, genre, price)

- *author → gender*
- *booktitle → genre, price*

Splitting **Books** using the FD *author → gender*:

- **Author** (author, gender)

  FD: *author → gender*  in BCNF!
- **Books2** (authos, booktitle, genre, price)

  FD: *booktitle → genre, price*  not in BCNF!

# BCNF Example (3)

**Books** (author, gender, booktitle, genre, price)

- $author \rightarrow gender$
- $booktitle \rightarrow genre, price$

Splitting **Books** using the FD $author \rightarrow gender$:
- **Author** (author, gender)
  FD: $author \rightarrow gender$  in BCNF!

- Splitting **Books2** (author, booktitle, genre, price):
  - **BookInfo** (booktitle, genre, price)
    FD: $booktitle \rightarrow genre, price$  in BCNF!
  - **BookAuthor** (author, booktitle) in BCNF!

# Third Normal Form (3NF)

# 3NF Definition

A relation **R** is in **3NF** if whenever $X \longrightarrow A$, one of the following is true:

- $A \in X$ (trivial FD)

- $X$ is a superkey

- $A$ is part of some key of **R** (prime attribute)

BCNF implies 3NF

# 3NF Cont'd

- Example: **R**(A, B, C) with $A, B \longrightarrow C$ and $C \longrightarrow A$
    - is in 3NF. Why?
    - is not in BCNF. Why?


- Compromise used when BCNF not achievable: *aim for BCNF and settle for 3NF*

- Lossless-join and dependency preserving decomposition into a collection of 3NF relations is always possible!

# 3NF Algorithm

1. Apply the algorithm for BCNF decomposition until all relations are in 3NF (we can stop earlier than BCNF)

2. Compute a minimal basis $F'$ of $F$

3. For each non-preserved FD $X \longrightarrow A$ in $F'$, add a new relation R(X, A)

# 3NF Example (1)

Start with relation **R** (A, B, C, D) with FDs:

- $A \longrightarrow D$
- $A, B \longrightarrow C$
- $A, D \longrightarrow C$
- $B \longrightarrow C$
- $D \longrightarrow A, B$

**Step 1**: find a BCNF decomposition

- **R1** (B, C)
- **R2** (A, B, D)

# 3NF Example (2)

Start with relation **R** (A, B, C, D) with FDs:

- $A \longrightarrow D$
- $A, B \longrightarrow C$
- $A, D \longrightarrow C$
- $B \longrightarrow C$
- $D \longrightarrow A, B$

**Step 2**: compute a minimal basis of the original set of FDs:

- $A \longrightarrow D$
- $B \longrightarrow C$
- $D \longrightarrow A$
- $D \longrightarrow B$

# 3NF Example (3)

Start with relation **R** (A, B, C, D) with FDs:

- $A \longrightarrow D$
- $A, B \longrightarrow C$
- $A, D \longrightarrow C$
- $B \longrightarrow C$
- $D \longrightarrow A, B$

**Step 3**: add a new relation for any FD in the basis that is not satisfied:

- all the dependencies in $F'$ are satisfied!
- the resulting decomposition **R1**, **R2** is also BCNF!

# Is Normalization Always Good?

- Example: suppose A and B are always used together, but normalization says they should be in different tables
  - decomposition might produce unacceptable performance loss
- Example: data warehouses
  - huge historical DBs, rarely updated after creation
  - joins expensive or impractical

# RECAP

- Bad schemas lead to redundancy
- To "correct" bad schemas: decompose relations
  - lossless-join
  - dependency preserving
- Desired normal forms
  - **BCNF**: only superkey FDs
  - **3NF**: superkey FDs + dependencies with prime attributes on the RHS