

In the Lecture Series Introduction to Database Systems

## Normalization

Presented by Stéphane Bressan

*Introduction to Database Systems*

### Learning Objectives

- Understand the rationale (anomalies) and definition of the main **normal forms** based on functional dependencies (**2NF**, **3NF** and **BCNF**)
- Be able to **decompose** (or **synthesize**) a schema into a **dependency preserving BCNF or 3NF**.

*lossless &*

### Anomalies

- Redundant storage
- Update anomalies
- Insertion anomalies
- Deletion anomalies

### Anomalies: Example

Assume that the position determines the salary:  
company      position → salary

eNumber	firstName	lastName	address	depart- ment	Position	salary
1XLJ3	Dewi	Srijaya	12a Jln Lempeng	Toys	Clerk	2000
4W3E	Isabel	Leong	10 Outram Park	Sports	Trainee	1200
3X0E	John	Smith	107 Clementi Rd	Toys	Clerk	2000
5SD2	Axel	Bayer	55 Cuscaden Rd	Sports	Trainee	1200
6RG5	Winnie	Lee	10 West Coast Rd	Sports	Manager	2500
755Y	Sylvia	Tok	22 East Coast Lane	Toys	Manager	2600
2SD3	Eric	Wei	100 Jurong drive	Toys	Assistant manager	2200
?	?	?	?	?	Security guard	1500

key

Insertion anomaly

Potential deletion anomaly

Update anomaly

Redundant storage

### Decomposition

- The **decomposition** of a **relation scheme R** is a **set** of **relation scheme R<sub>i</sub>** such that:

$$\bigcup_i R_i = R$$

*Do not lose any of the columns*  
*Union of all columns in all*  
*tables (Fragments)*  
*= Col.s in original table*

### Normalization: Example

#### employee

eNumber	firstName	lastName	address	depart- ment	Position
1XLJ3	Dewi	Srijaya	12a Jln Lempeng	Toys	Clerk
4W3E	Isabel	Leong	10 Outram Park	Sports	Trainee
3X0E	John	Smith	107 Clementi Rd	Toys	Clerk
5SD2	Axel	Bayer	55 Cuscaden Rd	Sports	Trainee
6RG5	Winnie	Lee	10 West Coast Rd	Sports	Manager
755Y	Sylvia	Tok	22 East Coast Lane	Toys	Manager
2SD3	Eric	Wei	100 Jurong drive	Toys	Assistant manager

key

- Redundant storage? ☐ NO
- Update anomaly? ☐ NO
- Deletion anomaly? ☐ NO
- Insertion anomaly? ☐ NO

#### salary

Position	salary
Clerk	2000
Trainee	1200
Manager	2500
Assistant manager	2200
Security guard	1500

key

### Lossless Decomposition: Example

- The decomposition is **lossless** if we can **recover the initial table**:

**SELECT** first\_name, last\_name, address,  
department, T2.position, salary  
**FROM** T2, T3  
**WHERE** T2.position = T3.position

keep all columns &  
preserve all data

I stitch  
back  
all  
data

Introduction to Database Systems

### Lossless Decomposition: Counter Example

Flight Number	Departure time	Arrival time	Origin	Destination
SG12	12h00	13h00+	SIN	CDG
TG414	15h50	16h30	SIN	JKT
TG415	12h00	14h20	BKK	SIN

Flight Number is the key

Original Table

Introduction to Database Systems

### Lossless Decomposition: Counter Example

Flight Number	Departure time	Arrival time	Origin	Destination
SG12	12h00	13h00+	SIN	CDG
TG414	15h50	16h30	SIN	JKT
TG415	12h00	14h20	BKK	SIN

Flight Number	Departure time	Origin
SG12	12h00	SIN
TG414	15h50	SIN
TG415	12h00	BKK

losing  
functional  
dependency

Departure time	Arrival time	Destination
12h00	13h00+	CDG
15h50	16h30	JKT
12h00	14h20	SIN

demo which departure time for what flight

Introduction to Database Systems

### Lossless, Dependency Preserving

- The **decomposition is lossy**
- And we **lost functional dependencies**:

$\{ \text{Flight Number} \} \rightarrow \{ \text{Arrival time, Destination} \}$

Non dependency preserving

Introduction to Database Systems

### Decomposition: Example

- Consider the relation scheme  $\{C, S, J, D, P, Q, V\}$
- with FDs
  - $\{C\} \rightarrow \{S, J, D, P, Q, V\}$
  - $\{J, P\} \rightarrow \{C\}$
  - $\{S, D\} \rightarrow \{P\}$

(Contracts (contractid, supplierid, projectid, deptid, partid, qty on line))

Introduction to Database Systems

### Decomposition: Example

- Consider the decomposition into  $\{C, S, J, D, Q, V\}$ ,  $\{S, D, P\}$
- The decomposition is **lossless**
  - We can **recover the initial relation** thanks to **S, D**
- However the functional dependency  $\{J, P\} \rightarrow \{C\}$  **is lost** across the two relations...
- This decomposition is **not dependency preserving**

P there co.  
 $\{S, D\} \rightarrow \{P\}$   
on both sides,  
tool to bring  
everything together

Introduction to Database Systems

### Dependency Preserving Decomposition

- The **decomposition** of a **relation scheme**
  - R** with FDs **F**
    - Is a set of relation schemes **R<sub>i</sub>** with FDs **F<sub>i</sub>**
- The decomposition is **dependency preserving** if and only if
 
$$(\cup_i F_i)^+ = F^+$$

Introduction to Database Systems

### Too Much Decomposition

- It might be **tempting** to **decompose to the extreme**
- Evaluation of queries** may be **inefficient** since it will involve **combining several relations**

↓ join causes inefficiency

Introduction to Database Systems

### Too Much Decomposition: Example

Flight Number	Departure time	Origin
SG12	12h00	SIN
TG414	15h50	SIN
TG415	12h00	BKK

Flight Number	Arrival time	Destination
SG12	13h00+	CDG
TG414	16h30	JKT
TG415	14h20	SIN

can actually decompose into 4 tables

Introduction to Database Systems

### Looking for a "Good" Design

- The designer needs guidelines:
  - Normalization theory**
    - Minimal redundancy and no anomalies
    - Lossless decompositions
    - Dependency preserving decompositions
- But ultimately the designer needs to look at the **workload** (the **queries** and their **efficiency requirement**)

Introduction to Database Systems

### Second Normal Form (2NF)

- R** is a **relation schema**, with the set **F** of **FDs**
- R** is in **2NF** if and only if
  - For all **X: X ⊆ R**
  - and, for all **A ∈ R**
  - such that there **exists** a **FD: X → {A}** in **F+**
- Then
  - A ∈ X** (the **FD is trivial**), or
  - X is not a proper subset** of a candidate key for **R**, or
  - A is part of some candidate key for R** (**A is called a prime attribute**)



Not a proper subset

Candidate key: minimal superkey

Introduction to Database Systems

### Second Normal Form (2NF)

- Consider the relation scheme **{A,B,C,D}** with the FDs:
  - {A,B} → {C,D}** and **{A,B} → {A,B,C,D}**
  - {A} → {D}**
- {A,B}** is a **primary key** (it is not a proper subset)
- {A}** is a **proper subset** of a **primary key** (Violation)
- {D}** is **not part of some candidate key** (FD is not a prime attribute)
- This scheme is **not in 2NF**

Introduction to Database Systems

## Second Normal Form (2NF)

Decomposition ✓  
lossless ✓  
A on both sides ✓  
No dependency loss ✓

- E.g., Relation scheme  $\{A, D\}$  with FDs
  - $\{A\} \rightarrow \{D\}$  *A is a candidate key*
  - $\{A\}$  is the primary key (i.e., not proper subset)
  - The scheme is in 2NF
- E.g., Relation scheme  $\{A, B, C\}$  with FDs
  - $\{A, B\} \rightarrow \{C\}$  *{A, B} candidate key*
  - $\{A, B\}$  is the primary key
  - The scheme is in 2NF

Introduction to Database Systems

## Third Normal Form (3NF)

- R is a relation schema, with the set F of FDs
- R is in 3NF if and only if
  - For all  $X: X \subset R$
  - And, for all  $A \in R$
  - such that there exists a FD:  $X \rightarrow \{A\}$  in F+
- Then
  - $A \in X$  (trivial FD), or
  - X is a superkey for R, or
  - A is part of some candidate key for R



3NF is always in 2NF  
3NF stricter cos of 2nd cond.

## Third Normal Form (3NF)

- Rationale:
  - If X is not a superkey for R then this dependency is partial and may cause redundancy
  - If A is not part of some candidate key for R then there is a transitive dependency: candidate key  $\rightarrow X \rightarrow A$
  - Insertion/deletion/update anomalies

Introduction to Database Systems

## Boyce-Codd Normal Form (BCNF)

- R is a relation schema, with the set F of FDs
- R is in BCNF if and only if
  - For all  $X: X \subset R$
  - And, for all  $A \in R$
  - such that there exists a FD:  $X \rightarrow \{A\}$  in F+
- Then
  - $A \in X$ , or
  - X is a superkey for R

No last condition *Stricter, removing a possibility*  
BCNF always in 3NF

Introduction to Database Systems

## BCNF $\subset$ 3NF $\subset$ 2NF $\subset$ 1NF \*

- 2NF:
  - $A \in X$ , or
  - X is not a proper subset of a candidate key for R, or
  - A is part of some candidate key for R
- 3NF:
  - $A \in X$ , or
  - X is a superkey for R, or
  - A is part of some candidate key for R
- BCNF:
  - $A \in X$ , or
  - X is a superkey for R *No last condition*

Introduction to Database Systems

All LHS must give ALL columns in F+ !!

## Decomposition into BCNF

Let S be the initial set of schemes with FDs F  
**Until** all relation schemes in S are in BCNF  
**for each** R in S  
 if FD  $X \rightarrow Y$  in F+ violates BCNF for R  
 ( $X \rightarrow R$  not hold, not superkey  
 $X \cap Y = \emptyset$ , not trivial)  $\rightarrow$  Violate conditions  
 then  
 use  $X \rightarrow X+$   
 let S be  $(S - \{R\}) \cup \{(R-X+) \cup X, X+\}$   
**endfor**  
**enduntil**

Introduction to Database Systems

Decomposition into BCNF

$X \rightarrow Y$   
 $X = \{c, d\}$   
 $X^+ = \{c, d, e, f\}$

$\{c, d\} \rightarrow \{e, f\}$  not trivial  
 $\{c, d\}$  not superkey

Find one culprit & break it into 2

a	b	c	d	e	f

a	b	c	d

c	d	e	f

*Introduction to Database Systems*

Decomposition into BCNF

- The **different possible orders\*** in which we may consider the dependencies violating BCNF in the algorithm application may **lead to different decompositions**

\*orders in which we consider the constraints violating the BCNF condition

*There may be problems which may lose functional dependencies. May get sth non-dependent preserving.*

*Introduction to Database Systems*

Remark: Projecting FDs

- If  $S$  is a fragment after decomposition of a relation  $R$  with FDs  $F$
- The set of projected FDs on  $S$  is the set  $G$  of FDs
  - If  $X \rightarrow Y$  is in  $G$ 
    - Then  $X$  and  $Y$  are subsets of  $S$
    - $X \rightarrow Y$  is in  $F^+$
  - If  $X \rightarrow Y$  is in  $F^+$  and  $X$  and  $Y$  are subsets of  $S$ 
    - Then  $X \rightarrow Y$  is in  $G$

*Introduction to Database Systems*

Decomposition into BCNF

- Let us consider the relation scheme  $R = \{A, B, C, D, E\}$  and the FDs:
  - $\{A\} \rightarrow \{B\}$
  - $\{A\} \rightarrow \{E\}$
  - $\{C\} \rightarrow \{D\}$

*Non-trivial, non-superkey*

$\{A, C\} \rightarrow \{A, B, C, D, E\}$
- Candidate key:  $\{A, C\}$ 

*BCNF Minimal cover*  
 $\{A\} \rightarrow \{B, E\}$   
 $\{C\} \rightarrow \{D\}$

*Decomp:*  
 $R_1: \{A, B, E\}$   
 $R_2: \{C, D\}$   
 $R_3: \{A, C\}$

*3NF algo: Calc minimal cover.*  
 $\hookrightarrow$  This case, already a minimal cover.

*Decomposition:*  $R_1 = \{A, B\}$   
 $R_2 = \{A, E\}$   
 $R_3 = \{C, D\}$   
 $R_4 = \{A, C\}$

*If 3NF synthesis, usually can find BCNF (NOT ALWAYS)*

*Introduction to Database Systems*

Decomposition into BCNF

- $R$  is not in BCNF
- Because (for instance):
  - $\{A\} \rightarrow \{B\}$  holds
  - It is not trivial
  - $\{A\}$  is not a superkey

*Introduction to Database Systems*

Decomposition into BCNF

- Pick  $\{A\} \rightarrow \{B\}$  for decomposition
- Expand into  $\{A\} \rightarrow \{B, E\}$
- $\{A, B, C, D, E\}$  becomes
  - $\{A, C, D\}$  and  $\{A, B, E\}$
- With FDs: (we need projected FDs)
  - $\{A\} \rightarrow \{B\}$ , on  $\{A, B, E\}$  ✓
  - $\{A\} \rightarrow \{E\}$ , on  $\{A, B, E\}$  ✓
  - $\{C\} \rightarrow \{D\}$ , on  $\{A, C, D\}$  ✓

*$\{A\}$  is superkey BCNF Violation*

*Introduction to Database Systems*

### Decomposition into BCNF

- Pick  $\{C\} \rightarrow \{D\}$  for decomposition
- Expand into  $\{C\} \rightarrow \{D\}$
- $\{A, C, D\}$  and  $\{A, B, E\}$  becomes
  - $\{A, C\}$ ,  $\{C, D\}$  and  $\{A, B, E\}$
- With FDs: (we need projected FDs)
  - $\{A\} \rightarrow \{B\}$ , on  $\{A, B, E\}$
  - $\{A\} \rightarrow \{E\}$ , on  $\{A, B, E\}$
  - $\{C\} \rightarrow \{D\}$ , on  $\{C, D\}$

*Introduction to Database Systems*

### Decomposition into BCNF

- Finally the BCNF decomposition of  $R=\{A, B, C, D, E\}$  with the FDs:
  - $\{A\} \rightarrow \{B\}$ ,
  - $\{A\} \rightarrow \{E\}$ ,
  - $\{C\} \rightarrow \{D\}$

- is:  $\{A, C\}$ ,  $\{C, D\}$  and  $\{A, B, E\}$

↳ Dependency preserving

*Introduction to Database Systems*

### Decomposition into BCNF

- BCNF decomposition may not be dependency preserving
- Example:  $\{A, B, C\}$  with FDs
  - $\{A, B\} \rightarrow \{C\}$
  - $\{C\} \rightarrow \{A\}$

Decompose to  $\{A, B\}$ ,  $\{C\}$ .
- The first FD is not preserved.

Divide & Conquer Algo

*Introduction to Database Systems*

### Decomposition into 3NF

Let S be the initial set of relation scheme R with FDs F

for each  $X \subset R$  such that  $X \rightarrow \{A_i\}$  is in the minimal cover of F  
if no scheme contains  $X \cup \{A_i\}$   
then create relation with scheme  $X \cup \{A_i\}$   
endfor

if no scheme contains a key for R  
then create a relation with scheme with key for R

*Introduction to Database Systems*

### Dependency preserving

- The 3NF synthesis algorithm always finds a lossless dependency preserving decomposition

*Introduction to Database Systems*

### Credits

The content of this lecture is based on chapter 8 of the book "Introduction to database Systems" By S. Bressan and B. Catania, McGraw Hill publisher

Clipart and media are licensed from Microsoft Office Online Clipart and Media

Copyright © 2013 by Stéphane Bressan



*Introduction to Database Systems*