

Databases

TDA357/DIT621– LP3 2023

Lecture 7

Ana Bove

(much of the material is based on material from
both Thomas Hallgren and Jonas Duregård)

February 2nd 2023

Recall Last Lecture

- Functional dependencies:
 - $X \rightarrow A$ if any two rows that agree on the values in X also agree on the values in A ;
 - Properties: reflexivity, transitivity and augmentation;
 - Closures: $X^+ = \{a \mid X \rightarrow a\}$;
 - Superkeys ($S \subseteq X^+$) and keys (minimal superkey);
 - Minimal basis/cover;
- Given $R(S)$, the FD $X \rightarrow A$ is a BCNF-violation if X is not a superkey ($X \subseteq S \subseteq X^+$);
- Given a (minimal basis) set of FD, $R(S)$ in BCNF if for all non-trivial $X \rightarrow A$, X is a superkey ($X \subseteq S \subseteq X^+$) (alt. if there is no BCNF-violation);
- Normalisation algorithm: given $R(S)$ and a BCNF-violation $X \rightarrow A$, we split $R(S)$ into $R_1(X^+)$ and $R_2(X \cup (S - X^+))$ and recursively normalise them (R and S become now irrelevant).

Overview of Today's Lecture

- FD examples;
- BCNF decomposition and dependency preservation;
- Multivalued dependencies (MVD);
- 4NF and its normalisation algorithm;
- MVD examples.

Example: FD and BCNF-violation

Consider a relation schema $R(a, b, c, d, e, f)$ with the following FD. Identify 3 different BCNF-violations and compute the closure of the left-hand side of the dependency.

$a \rightarrow b$
 $ac \rightarrow d$
 $ac \rightarrow e$
 $ac \rightarrow f$
 $d \rightarrow a$
 $d \rightarrow c$
 $af \rightarrow e$
 $e \rightarrow a$

$$a \rightarrow b: \{a\}^+ = \{a, b\}.$$

$$af \rightarrow e: \{a, f\}^+ = \{a, b, e, f\}.$$

$$e \rightarrow a: \{e\}^+ = \{a, b, e\}.$$

For the rest we have that $\{a, c\}^+ = \{d\}^+ = \{a, b, c, d, e, f\}$.

Example: FD and Data

Consider the schema $R(a, b, c, d)$
with the following FD:

$a \rightarrow b$
 $b \rightarrow a$
 $c \rightarrow d$

Construct a table with these FD
and no other FD with only one at-
tribute on the left side.

Give 2 possible sets of keys.

a	b	c	d
1	7	2	2
1	7	3	3
4	4	2	2
5	6	5	2

Keys: $\{a, c\}$ or $\{b, c\}$

The solution could look like this:

a	b	c	d
1	1	2	2
1	1	2	3

This solution is however incorrect,
since it for example violates $c \rightarrow d$.
Also, $a \rightarrow c$ is a FD but shouldn't.

The first two lines make sure that there are no
dependencies from a and/or b to c or d .

The third line does the same in the other direction.

The last makes sure that $d \rightarrow c$ does not hold.

Observe that a and b do not need to be equal!

BCNF Decomposition and Dependency Preservation

Given the FD:

$$\begin{array}{l} x \rightarrow z \\ yz \rightarrow w \end{array}$$

normalise

$$R(x, y, z, w)$$

- We obtain the derived FD $xy \rightarrow w$:
 - from $x \rightarrow z$ we get $xy \rightarrow z$ by augmentation;
 - we have $y \rightarrow y$ by reflexivity and then $xy \rightarrow y$ by augmentation;
 - $xy \rightarrow z$ and $xy \rightarrow y$ give us $xy \rightarrow yz$;
 - by transitivity of $xy \rightarrow yz$ and $yz \rightarrow w$ we get $xy \rightarrow w$.
- Using $x \rightarrow z$ to decompose and the FD (including derived ones) to determine the keys, we get:
 - $R_1(\underline{x}, z)$ in BCNF
 - $R_2(\underline{x}, \underline{y}, w)$ also in BCNF
- Note that $yz \rightarrow w$ is not guarantee!

Multivalued Dependencies (MVD)

Consider the following example with courses, books, authors and teachers such that:

- Databases has a book and two teachers;
- Automata has a teacher and two books.

course	book	author	teacher
Databases	DTCB	Ullman	Jonas
Databases	DTCB	Ullman	Ana
Databases	DTCB	Widom	Jonas
Databases	DTCB	Widom	Ana
Automata	Finite Automata	Carroll	Nisse
Automata	Automata Theory	Hopcroft	Nisse

This relation has no non-trivial FD...

... however, there is risk for update and deletion anomalies!

Multivalued Dependencies (Cont.)

course	book	author	teacher
Databases	DTCB	Ullman	Jonas
Databases	DTCB	Ullman	Ana
Databases	DTCB	Widom	Jonas
Databases	DTCB	Widom	Ana
Automata	Finite Automata	Carroll	Nisse
Automata	Automata Theory	Hopcroft	Nisse



- Course determines the set of teachers, which are independent of the other values (books and authors);
- Course determines the set of books and authors, which are independent of the other values (teachers);
- There might be *multiple values* of teachers/books and authors related to a course.



course \twoheadrightarrow teacher
course \twoheadrightarrow book author

Multivalued Dependencies (MVD)

Informally, a MVD $X \twoheadrightarrow Y$ holds in a relation for $R(a_1, \dots, a_n)$ if whenever *there are tuples* $\langle \bar{x}, \bar{y}_1, \bar{z}_1 \rangle$ and $\langle \bar{x}, \bar{y}_2, \bar{z}_2 \rangle$ then *there are also tuples* $\langle \bar{x}, \bar{y}_1, \bar{z}_2 \rangle$ and $\langle \bar{x}, \bar{y}_2, \bar{z}_1 \rangle$ with

- \bar{x} are the values for the attributes in X ;
- \bar{y}_1 and \bar{y}_2 are the values for the attributes in Y ;
- and \bar{z}_1 and \bar{z}_2 are the values for the rest of the attributes (those in $Z = \{a_1, \dots, a_n\} - X - Y$).

That is, we have all possible combination of values for Y (\bar{y}_1 and \bar{y}_2) and Z (\bar{z}_1 and \bar{z}_2) for those rows with the same value in X !

MVD are not easy to figure out nor to check!

Multivalued Dependencies (MVD)

Definition: A *MVD* $X \twoheadrightarrow Y$ *holds* in a relation for $R(a_1, \dots, a_n)$ if for any tuples t_1 and t_2 such that $t_1.X = t_2.X$, there exists tuples t_3 and t_4 such that

- $t_1.X = t_2.X = t_3.X = t_4.X$;
- $t_1.Y = t_3.Y$ and $t_2.Y = t_4.Y$;
- and for $Z = \{a_1, \dots, a_n\} - X - Y$ (the rest of the attributes), then $t_1.Z = t_4.Z$ and $t_2.Z = t_3.Z$.

t_3 and t_4 could actually be the same as t_1 and t_2 for some of the choices of t_1 and t_2 !

Multivalued Dependencies (MVD): Back to the Example

row	course	book	author	teacher
r_1	Databases	DTCB	Ullman	Jonas
r_2	Databases	DTCB	Ullman	Ana
r_3	Databases	DTCB	Widom	Jonas
r_4	Databases	DTCB	Widom	Ana
r_5	Automata	Finite Automata	Carroll	Nisse
r_6	Automata	Automata Theory	Hopcroft	Nisse

course \twoheadrightarrow teacher: $X = \{\text{course}\}$, $Y = \{\text{teacher}\}$, $Z = \{\text{book}, \text{author}\}$ (the rest)

$$r_1.X = r_2.X = r_3.X = r_4.X = \text{Databases}$$

$$r_1.Y = r_3.Y = \text{Jonas and } r_2.Y = r_4.Y = \text{Ana}$$

$$r_1.Z = r_2.Z = \text{DTCB+Ullman and } r_3.Z = r_4.Z = \text{DTCB+Widmon}$$

course \twoheadrightarrow book author: $X = \{\text{course}\}$, $Y = \{\text{book}, \text{author}\}$, $Z = \{\text{teacher}\}$ (the rest)

$$r_1.X = r_2.X = r_3.X = r_4.X = \text{Databases}$$

$$r_1.Y = r_2.Y = \text{DTCB+Ullman and } r_3.Y = r_4.Y = \text{DTCB+Widmon}$$

$$r_1.Z = r_3.Z = \text{Jonas and } r_2.Z = r_4.Z = \text{Ana}$$

For the rows with same value in X , the cartesian product with values of Y and Z should be in the table!

Properties of Multivalued Dependencies

Let $R(a_1, \dots, a_n)$ be a relational schema, $S = \{a_1, \dots, a_n\}$ the set of attributes of R , and $X, Y, V, W \subseteq S$.

Replication: if $X \rightarrow Y$ then $X \twoheadrightarrow Y$;

Complementation: If $X \twoheadrightarrow Y$ then $X \twoheadrightarrow S - Y$;

Transitivity: If $X \twoheadrightarrow Y$ and $Y \twoheadrightarrow W$ then $X \twoheadrightarrow W - Y$;

Augmentation: If $X \twoheadrightarrow Y$ and $V \subseteq W$ then $XW \twoheadrightarrow YV$;

Trivial: $X \twoheadrightarrow Y$ is trivial if $Y \subseteq X$ or if $S \subseteq X \cup Y$.

Observe the difference in the definition of a trivial dependency!

4th Normal Form (4NF)

Definition: A relational schema $R(S)$ with set of attributes $S = \{a_1, \dots, a_n\}$ is in *4th normal form (4NF)* if

- The relational schema is in BCNF;
- For all non-trivial MVD $X \twoheadrightarrow Y$, X is a superkey (that is, it determines –following the FD– all the attributes in S).

If X is not a superkey, then $X \twoheadrightarrow Y$ is a *4NF-violation*.

4NF Normalisation Algorithm

To *normalise* a relational schema $R(S)$ with $S = \{a_1, \dots, a_n\}$:

- Find a 4NF-violation: that is, a non-trivial MVD $X \twoheadrightarrow Y$ ($Y \not\subseteq X$ or $S \not\subseteq X \cup Y$) such that X is not a superkey ($S \not\subseteq X^+ = \{a \mid X \rightarrow a\}$);
- If there is no such MVD then R is already in 4NF;
- Otherwise decompose $R(S)$ into $R_1(X \cup Y)$ and $R_2(S - Y)$ and normalise them both.

Note: Again, R is not of interest anymore since it has been replaced by R_1 and R_2 and the set S of attributes is replaced by $X \cup Y$ and $S - Y$ respectively.

Decomposition of Data

R : Courses

course	book	author	teacher
Databases	DTCB	Ullman	Jonas
Databases	DTCB	Ullman	Ana
Databases	DTCB	Widom	Jonas
Databases	DTCB	Widom	Ana
Automata	Finite Automata	Carroll	Nisse
Automata	Automata Theory	Hopcroft	Nisse



course \twoheadrightarrow teacher



R_1 : CourseTeachers

course	teacher
Databases	Jonas
Databases	Ana
Automata	Nisse

R_2 : CourseBooks

course	book	author
Databases	DTCB	Ullman
Databases	DTCB	Widom
Automata	Finite Automata	Carroll
Automata	Automata Theory	Hopcroft

Lossless Join: All Data Back!

R_1 : CourseTeachers

course	teacher
Databases	Jonas
Databases	Ana
Automata	Nisse

R_2 : CourseBooks

course	book	author
Databases	DTCB	Ullman
Databases	DTCB	Widom
Automata	Finite Automata	Carroll
Automata	Automata Theory	Hopcroft

Query: R_1 NATURAL JOIN R_2

R : Courses

course	book	author	teacher
Databases	DTCB	Ullman	Jonas
Databases	DTCB	Ullman	Ana
Databases	DTCB	Widom	Jonas
Databases	DTCB	Widom	Ana
Automata	Finite Automata	Carroll	Nisse
Automata	Automata Theory	Hopcroft	Nisse

What about the Primary Keys?

R_1 : CourseTeachers

course	teacher
Databases	Jonas
Databases	Ana
Automata	Nisse

R_2 : CourseBooks

course	book	author
Databases	DTCB	Ullman
Databases	DTCB	Widom
Automata	Finite Automata	Carroll
Automata	Automata Theory	Hopcroft



$R_1(\text{course}, \underline{\text{teacher}})$

$R_2(\text{course}, \underline{\text{book}}, \underline{\text{author}})$

MVD Example

Identify a non-trivial MVD that holds in the data and is a 4NF-violation, and provide a schema in 4NF based on that MVD.

code	examDate	student	classification
TDA357	20200115	Emil	CS
TDA357	20200115	Emilia	CS
TDA357	20200320	Emil	CS
TDA357	20200115	Emil	DB
TDA357	20200115	Emilia	DB
TDA357	20200320	Emil	DB
XYZ123	20200115	Emil	CS

$\text{code} \twoheadrightarrow \text{classification}$ (or equivalently: $\text{code} \twoheadrightarrow \text{examDate}, \text{student}$)

$R_1 (\underline{\text{code}}, \underline{\text{classification}})$

$R_2 (\underline{\text{code}}, \underline{\text{examDate}}, \underline{\text{student}})$

MVD Example (from exam)

Consider the relation schema $R(g,h,i,j,k)$ and the MVD $g \ h \twoheadrightarrow i$.

- ① Complete the following table so that the dependency is valid

g	h	i	j	k
0	0	1	2	2
0	0	1	3	3
0	0	4	2	2
5	5	1	2	2
5	5	1	3	3

- ② Convert the schema into a relation schema in 4NF.

MVD Example (from exam, cont.)

1

g	h	i	j	k
0	0	1	2	2
0	0	1	3	3
0	0	4	2	2
0	0	4	3	3
5	5	1	2	2
5	5	1	3	3

- 2 Using this MVD we obtain the schemas $R1(\underline{g}, \underline{h}, \underline{i})$ and $R2(\underline{g}, \underline{h}, \underline{j}, \underline{k})$.

Database Design Workflow

- ER-modelling and FD are two complementary (but compatible) methods to find a good relational schema;
- FD can find certain things that ER-diagram cannot, and the other way around;

FD are particularly useful to find secondary keys (**UNIQUE** constraints);

- A good database design workflow is:
 - Start by creating the ER-model from the domain description;
 - Obtain a relational schema from the ER-model;
 - Go back to the domain description and identify additional FD;
 - Check if they hold in the ER-schema and add secondary keys when appropriate;
 - Apply BCNF decomposition when needed;
 - Identify MVD and apply 4NF decomposition when needed.

Other Normal Forms

1NF (1970): Only atomic values in the tables;

2NF (1971): 1NF + no partial dependencies,
(a FD $X \rightarrow Y$ is a partial dependency if $Z \rightarrow Y$ for some $Z \subseteq X$);

3NF (1971): 2NF + no transitive dependencies from the primary key to non-key attributes,
Schemas that result from ER-modelling are in 3NF;

BCNF a.k.a. 3^{1/2}NF (1974): for each non-trivial functional dependency $X \rightarrow Y$, X is a superkey;
(BCNF implies 2NF and also 3NF!);

4NF (1977): BCNF + no MVD violation;

5NF, 6NF, ...

Overview of Next Lecture

- Functions;
- Triggers:
 - On tables;
 - On views;
- Example.

Reading:

Book: chapter 7.5

Notes: chapter 7.4.4–7.4.6