Databases

TDA357/DIT621-LP3 2023

Lecture 6

Ana Bove

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

January 30th 2023

Recall Last Lecture

- Entity-Relationship modelling:
 - Entities and attributes;
 - Many-to-many relationships;
 - Many-to-exactly-one relationships;
 - Many-to-at-most-one relationships;
 - Multiway relationships;
 - Self-relationships;
 - Weak entities;
 - ISA relationships;
- (ER-example/exercise.)

• Functional dependencies:

- Reflexivity, transitivity and augmentation;
- Closures;
- Superkeys and keys;
- Minimal basis/cover;
- Boyce-Codd Normal form (BCNF) ..
- ... and its normalisation algorithm.

Functional Dependencies and Normal Forms

This is an alternative, and in some way complementary, approach to model a database.

We start with a domain description and end up with a database schema.



Relations, Relation Schemas and Tables

Recall: A relation *S* is a subset of the cartesian product of two or more sets T_1, T_2, \ldots, T_n :

 $S \subseteq T_1 \times T_2 \times \cdots \times T_n$

Given a relation schema $R(a_1, \ldots, a_n)$, consider the *domain/type* of each attribute $a_1 : T_1, \ldots, a_n : T_n$.

The *relation signature* of the relation R is the corresponding cartesian product $T_1 \times \cdots \times T_n$.

So, given a relation schema $R(a_1, \ldots, a_n)$ with signature $T_1 \times \cdots \times T_n$:

• A *table* for the schema $R(a_1, ..., a_n)$ is a subset of the cartesian product $T_1 \times \cdots \times T_n$;

• A *row* in the table is an element of the cartesian product $t \in T_1 \times \cdots \times T_n$.

Attribute Names vs Tuple Components

Recall: The elements of a cartesian product $T_1 \times \cdots \times T_n$ are tuples $t = \langle t_1, \ldots, t_n \rangle$.

The *i*th projection $\pi_i t$ gives us the element t_i (of type T_i).

If t is a row in the table for the schema $R(\ldots, a, \ldots)$ (containing attribute a), we will:

- Assume an *indexing* function *i* giving us the index/position of each attribute;
- Denote $t.a = t_{i(a)}$ the projection $\pi_{i(a)} t$;

• If
$$A = \{a_1, \ldots, a_n\}$$
 is a set of attributes in $R(\ldots)$, then $t.A = \langle t.a_1, \ldots, t.a_n \rangle$ is the *simultaneous projection*.

CREATE TABLE Countries (name TEXT ..., abbr CHAR(2) ..., area FLOAT ...); Relation schema: Countries (name, abbr, area)

 $\frac{\text{Relation signature:}}{\text{TEXT} \times \text{CHAR}(2) \times \text{FLOAT}}$

name	abbr	area
Denmark	DK	43094
Sweden	SE	449964
Norway	NO	324220

Let $t = \langle$ 'Denmark', 'DK', 43094 \rangle

If $A = \{abbr, area\}$ then $t.A = \langle 'DK', 43094 \rangle$

Functional Dependencies (FD)

Let $R(a_1, \ldots, a_n)$ be a relation schema, $S = \{a_1, \ldots, a_n\}$ the set of attributes of R, and $X, A \subseteq S$.

Definition: (*Functional Dependency*) X determines A, denoted $X \rightarrow A$, iff for all rows $t, u \in R(...)$, if t.X = u.X then t.A = u.A.

Example: Suppose we have $a b \rightarrow c$. What does this mean?

It means that if two rows agree on the values of the attributes a and b, then they should also agree on the values of the attributes c.

Hence, the values of *a* and *b* uniquely determine the values of *c*.

Capital vs small letters: We will use capital letters A, ..., X, Y, Z to denote sets of attributes and small letters a, b, c, ..., z to denote single attributes (unless otherwise stated).

 $a \rightarrow b c$: means that a determines both b and c:

$$a \rightarrow b c$$
 is the same as $a \rightarrow b$
 $a \rightarrow c$

 $a b \rightarrow c$: means that a and b together determine c:

$$a b \rightarrow c$$
 is NOT the same as

$$\begin{array}{c}
 a \to c \\
 b \to c
\end{array}$$

Properties of Functional Dependencies

Let $R(a_1, \ldots, a_n)$ be a relation schema, $S = \{a_1, \ldots, a_n\}$ the set of attributes of R, and $X, Y, Z, A \subseteq S$.

Transitivity: If $X \to Y$ and $Y \to Z$ then $X \to Z$;

Augmentation: If $X \to A$ and a is an attribute, then $X a \to A$ and $a X \to A$;

Reflexivity: (trivial dependency) $X \rightarrow X$;

Reflexivity + augmentation (trivial dependency) If $Y \subseteq X$ then $X \to Y$.

Closure: $X^+ = \{a \mid X \to a\}$, the set of all attributes determined by X;

Superkey: X such that $S \subseteq X^+$; (X is a superkey if X^+ includes all the attributes in R(...))

(Minimal) Key: Minimal superkey (and good candidate for primary key!); removing an attribute to the set will make it a non-superkey.

Example: Deriving FD

Given the FD:

$$egin{array}{c} x
ightarrow y \ z
ightarrow w \ y w
ightarrow q \end{array}$$

we can derive the FD:

$$xz \rightarrow q$$

- $x z \rightarrow y$: from $x \rightarrow y$ by augmentation;
- $x z \rightarrow w$: from $z \rightarrow w$ by augmentation;
- $x z \rightarrow y w$: by merging left-hand side of FD (see words on notation);
- $x z \rightarrow q$: by transitivity of $x z \rightarrow y w$ and $y w \rightarrow q$.

Example: Computing the Closure

Given the FD:

$ \begin{array}{c} x \to y \\ z \to w \\ y w \to q \\ q \to x \\ r \to s \end{array} $	compute:	${x,z}^+$
--	----------	-----------

The closure gives us the non-trivial FD: (Recall: if $Y \subseteq X$ then $X \rightarrow Y$ trivial)

$$\left(\begin{array}{c}
x \, z \to y \\
x \, z \to w \\
x \, z \to q
\end{array}\right)$$

$$x z \rightarrow y w q$$

or

There are three ways we can use functional dependencies:

- Check if they hold for a specific data set;
- Check if a specific design/relational schema ensures the FD hold for all data set that follows the schema;
- Express desired properties a design/relational schema should have.

(This use is what makes FD a design tool, and what we will concentrate on in (most of) the rest of this lecture.)

Example: Which FD Hold for this Data?

code	name	day	time	room	seats
TMV028	Automata	Monday	10	HB1	108
TDA357	Databases	Monday	15	HC4	104
TMV028	Automata	Tuesday	13	HB1	108
TDA357	Databases	Wednesday	10	HC2	115
TDA357	Databases	Thursday	10	HC4	104

 $code \rightarrow name: YES$



day time room \rightarrow code: YES



Example: Which FD Hold for this Design?

 $\begin{array}{l} \mbox{Teachers (name, email)} \\ \mbox{Courses (code, cname, teacher)} \\ \mbox{teacher} \rightarrow \mbox{Teachers.name} \end{array}$

Does the relational schema guarantee:



Example: Deriving Closures, Keys and Superkeys (I)



Example: Deriving Closures, Keys and Superkeys (II)



Example: Deriving Closures, Keys and Superkeys (III)



Example: Deriving Closure, Keys and Superkeys (III, Cont.)

If country is the primary key, then the dependency

 $country \rightarrow currency$

will be trivially satisfied.

But how to ensure that

 $\mathsf{currency} \to \mathsf{value}$

is satisfied? This is a non-trivial FD and {currency, value} is not a superkey!

Note: FD help us identify a problematic schema!

Let F be a set of functional dependencies.

Definition: The *minimal basis or minimal cover* F^- is a simplified but equivalent set of functional dependencies such that:

- F^- contains no trivial dependencies (if $Y \subseteq X$ then $X \to Y$ trivial);
- No dependencies in *F*⁻ follow from other dependencies in *F*⁻ through transitivity or augmentations.

Example: Deriving Minimal Basis/Cover

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

Given the FD:

we can compute the minimal basis by removing:

- $a d \rightarrow d$: because it is trivial;
- $a d \rightarrow b$: because it can be computed from $a \rightarrow b$ by augmentation;
- $a d \rightarrow c$: because it can be computed from $a \rightarrow b$ and $b \rightarrow c$ by transitivity, and then augmentation.

Minimal basis/cover:

$$\left(\begin{array}{c} a \to b \\ b \to c \end{array}\right)$$

Deriving Minimal Basis/Cover

A way to see if a FD $X \rightarrow Y$ can be derived from the other FD is to check whether X^+ is a superkey when $X \rightarrow Y$ is not taken into account.

Example: Consider the FD:
$$a \rightarrow b \\ b c \rightarrow d \\ a c \rightarrow d$$

Is $a c \rightarrow d$ derived?

Let us compute
$$\{a, c\}^+$$
 from $\begin{array}{c} a \to b \\ b c \to c \end{array}$

$$\begin{bmatrix} \mathsf{a} \to \mathsf{b} \\ \mathsf{b} \, \mathsf{c} \to \mathsf{d} \end{bmatrix}$$

$$\{a,c\}^+ = \{a,b,c,d\}$$

That is, $\{a, c\}^+$ is a superkey and hence it should be possible to derived $a c \rightarrow d$ from the other FD.

Minimal basis/cover: is $F^- \subseteq F$?

Given a set F of FD, the minimal basis/cover F^- does not have to be a subset of F.

Consider the FD:

$$a c o b$$

 $a o c$

Minimal basis/cover:

$$egin{array}{c} a
ightarrow b \ a
ightarrow c \end{array}$$

since $a \rightarrow c$, then $ac \rightarrow b$ can actually be derived from $a \rightarrow b$ by augmentation. Otherwise $ac \rightarrow b$ might hold but not $a \rightarrow b$!

The minimal basis is not included in the original set of FD.

FD: Summary so Far

- A FD $X \rightarrow Y$ means that any rows that agree on X also agree on Y;
- We can extend a set of FD with additional derived FD using transitivity, augmentation, and reflexivity;
- Conversely, we can reduce a set of FD to its minimal basis/cover by removing all trivial and derived FD;
- The closure X⁺ is the set of all attributes that can be determined by X;
- A superkey is a set of attributes that determines all others in the relation;
- Keys are minimal superkeys;
- To find <u>a</u> key: start with all attributes (a trivial superkey) and remove attributes until it is a key (finding all keys is more work though).

- We work like this:
 - We start by collecting all the attributes in the domain and place them in one big relation schema R(a₁,..., a_n);
 - We collect the set F with all the FD;
 - (We compute the minimal basis/cover *F*⁻;)
 - We normalise R(...) using F (alt. F^-) to get a good design.
- Normalisation is a recursive procedure. To normalise R(...):
 - We check if R(...) is already in *normal form*;
 - If not, we decompose $R(\ldots)$ into $R_1(\ldots)$ and $R_2(\ldots)$ and normalise them.
- There are several *normal form* definitions ...
- ... and several *normalisation* algorithms (depending on the normal form definitions).

Definition: Given a relation schema $R(a_1, ..., a_n)$, the non-trivial FD $X \rightarrow Y$ (with $X \subseteq \{a_1, ..., a_n\}$) is a *BCNF-violation* if X is not a superkey ($\{a_1, ..., a_n\} \not\subseteq X^+$).

Definition: A relation schema $R(a_1, ..., a_n)$ is in *BCNF* if for each non-trivial FD $X \to Y$, X is a superkey $(X \subseteq \{a_1, ..., a_n\} \subseteq X^+)$.

That is, if there are no BCNF-violations.

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

- Find a BCNF-violation: that is, a non-trivial FD X → Y such that X is not a superkey (X ⊆ S ∉ X⁺);
- If there is no such FD then R is already in BCNF;
- Otherwise decompose R(S) into R₁(X⁺) and R₂(X ∪ (S − X⁺)) and normalise them both.

Note: R(S) is of no interest anymore and has been replaced by $R_1(S_1)$ and $R_2(S_2)!$

Example: BCNF Normalisation



Decompose using code \rightarrow name?

$$X = \{\text{code}\}, X^+ = \{\text{code}\}^+ = \{\text{code, name}\}$$
$$R_1(X^+) = \underbrace{R_1(\text{code, name})}_{(\text{in BCNF!})} \text{(in BCNF!)}$$
$$R_2(X \cup (S - X^+)) = R_2(\text{code, day, time, room, seats})$$

Decompose using room \rightarrow seats?

$$X = \{\text{room}\}, X^{+} = \{\text{room}\}^{+} = \{\text{room, seats}\}$$
$$R_{21}(X^{+}) = \boxed{R_{21}(\text{room, seats})} \text{ (in BCNF!)}$$
$$R_{22}(X \cup (S - X^{+})) = R_{22}(\text{code, day, time, room})$$

Example: BCNF Normalisation (Cont.)

Can R_{22} (code, day, time, room) be decomposed further?

We look at the remaining FD:

Decompose using day time room \rightarrow code? $X = \{ day, time, room \}$ $X^+ = \{ day, time, room \}^+ = \{ day, time, room, code, seats \}$

Decompose using day time code \rightarrow room?

 $X = \{$ day, time, code $\}$ $X^+ = \{$ day, time, code $\}^+ = \{$ day, time, room, code, name $\}$

So $[R_{22}(\text{code, day, time, room})]$ is in BCNF!

Both {day, time, room} and {day, time, code} are keys!

Identifying the Keys, Uniqueness Constrains and References



Relational Schema: R_1 (code, name) R_{21} (room, seats) R_{22} (code, day, time, room)

Keys can be determined using the FD and the normalisation algorithm. (A (minimal) X that determines all the attributes in the relation schema is a key!)

We get 2 possible solutions:

 $\begin{array}{l} R_1(\underline{\mathrm{code}},\,\mathrm{name})\\ R_{21}(\underline{\mathrm{room}},\,\mathrm{seats})\\ R_{22}(\underline{\mathrm{code}},\,\underline{\mathrm{day}},\,\underline{\mathrm{time}},\,\mathrm{room})\\ \mathrm{code} \to \overline{R_1}.\mathrm{code}\\ \mathrm{room} \to R_{21}.\mathrm{room}\\ \mathrm{Unique}\;(\mathrm{day},\,\mathrm{time},\,\mathrm{room}) \end{array}$

 $R_1(\underline{code}, name)$ R_{21} (room, seats) R_{22} (code, day, time, room) $code \rightarrow R_1.code$ room $\rightarrow R_{21}$.room Unique (day, time, code)

Note: We need to keep track of the splitting for the references!

BCNF Decomposition Avoids Data Redundancy (based on FD)

R: CourseBookings

code	name	day	time	room	seats
TMV028	Automata	Monday	10	HB1	108
TDA357	Databases	Monday	15	HC4	104
TMV028	Automata	Tuesday	13	HB1	108
TDA357	Databases	Wednesday	10	HC2	115
TDA357	Databases	Thursday	10	HC4	104



R₁: Courses



R₂₂: Bookings

ſ	code	name	room
	TMV028	Automata	HB1
	TDA357	Databases	HC2
- 1			

room	seats
HB1	108
HC2	115
HC4	104

code	day	time	room
TMV028	Monday	10	HB1
TDA357	Monday	15	HC4
TMV028	Tuesday	13	HB1
TDA357	Wednesday	10	HC2
TDA357	Thursday	10	HC4

Lossless Join: All Data Back!



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

Reading:

Book: chapter 3 Notes: chapters 4.1–4.3 and 5