

A “good” database schema should not lead to update anomalies.

Functional dependencies

A function f from S_1 to S_2 has the property

$$\text{if } x, y \in S_1 \text{ and } x = y, \text{ then } f(x) = f(y).$$

A generalization of keys to avoid design flaws violating the above rule.

Let X and Y be sets of attributes in R.

X (*functionally*) determines Y , $X \rightarrow Y$, iff $t_1[X] = t_2[X]$ implies $t_1[Y] = t_2[Y]$.

i.e., $f(t(X)) = t[Y]$

We also say $X \rightarrow Y$ is a *functional* dependency, and that Y is *functionally* dependent on X.

X is called the *left side*, Y the *right side* of the dependency.

Armstrong's axioms (1974)

Notation: If X and Y are sets of attributes, we write XY for their union.

e.g. $X = \{A, B\}$, $Y = \{B, C\}$, $XY = \{A, B, C\}$

F1 (Reflexivity) If $X \supseteq Y$ then $X \rightarrow Y$.

F2 (Augmentation) $\{X \rightarrow Y\} \models XZ \rightarrow YZ$.

F3 (Transitivity) $\{X \rightarrow Y, Y \rightarrow Z\} \models X \rightarrow Z$.

F4 (Additivity) $\{X \rightarrow Y, X \rightarrow Z\} \models X \rightarrow YZ$.

F5 (Projectivity) $\{X \rightarrow YZ\} \models X \rightarrow Y$.

F6 (Pseudotransitivity) $\{X \rightarrow Y, YZ \rightarrow W\} \models XZ \rightarrow W$.

To compute closure

$F = \{A \rightarrow B, BC \rightarrow D, A \rightarrow C\}$, compute $\{A\}^+$

1st scan of F:

$X^+ := \{A\}$

$X^+ := \{A, B\}$

$X^+ := \{A, B, C\}$

2nd scan of F:

$X^+ := \{A, B, C, D\}$

3rd scan of F: no change, therefore the algorithm terminates.

$\{A\}^+ := \{A, B, C, D\}$

To compute a Candidate key

Given a relational schema R and a set F of functional dependencies on R .

A key X of R must have the property that $X^+ = R$.

Algorithm to compute a candidate key

Step 1: Assign X a superkey in F .

Step 2: Iteratively remove attributes from X while retaining the property $X^+ = R$ till no reduction on X .

The remaining X is a key.

$$R = \{A, B, C, D\} \text{ and } F = \{ A \rightarrow B, BC \rightarrow D, A \rightarrow C \}$$

$X = \{A, B, C\}$ if the left hand side of F is a super key.

A cannot be removed because $\{BC\}^+ = \{B, C, D\} \neq R$

B can be removed because $\{AC\}^+ = \{A, B, C, D\} = R$
→ $X = \{A, C\}$

C can be further removed because $\{A\}^+ = \{A, B, C, D\}$
→ $X = \{A\}$

Normal Forms for Relational Databases

- criteria for a good database design (i.e., to resolve update anomalies)
- formalized by functional (or other) dependencies

Non-first normal form

Fac_Dept	Prof	Course Preferences	
		Course	Course_Dept
Comp Sci	Smith	353	Comp Sci
		379	Comp Sci
		221	Decision Sci
	Clark	353	Comp Sci
	351	Comp Sci	
	379	Comp Sci	
	456	Mathematics	
Chemistry	Turner	353	Comp Sci
		456	Mathematics
		272	Chemistry
Mathematics	Jameison	353	Comp Sci
		379	Comp Sci
		221	Decision Sci
		456	Mathematics
		469	Mathematics

First Normal Form (1NF)

This simply means that attribute values are atomic, and is part of the definition of the relational model.

Atomic: multivalued attributes, composite attributes, and their combinations are disallowed.

CRS_PREF			
Prof	Course	Fac_Dept	Crs_Dept
Smith	353	Comp Sci	Comp Sci
Smith	379	Comp Sci	Comp Sci
Smith	221	Comp Sci	Decision Sci
Clark	353	Comp Sci	Comp Sci
Clark	351	Comp Sci	Comp Sci
Clark	379	Comp Sci	Comp Sci
Clark	456	Comp Sci	Mathematics
Turner	353	Chemistry	Comp Sci
Turner	456	Chemistry	Mathematics
Turner	272	Chemistry	Chemistry
Jamieson	353	Mathematics	Comp Sci
Jamieson	379	Mathematics	Comp Sci
Jamieson	221	Mathematics	Decision Sci
Jamieson	456	Mathematics	Mathematics
Jamieson	469	Mathematics	Mathematics

Second Normal Form (2NF)

A *prime* attribute is one that is part of a candidate key. Other attributes are *non-prime*.

Definition: In an FD $X \rightarrow Y$, Y is *fully functionally dependent* on X if there is no $Z \subset X$ such that $Z \rightarrow Y$. Otherwise Y is *partially dependent* on X .

Proper Subset

Definition (*Second Normal Form*): A relation scheme is in second normal form (2NF) if all non-prime attributes are fully functionally dependent on the candidate keys.

A database scheme is in 2NF if all its relations are in 2NF.

Possible 2NF decomposition of the relation above is:

COURSE_PREF	
Prof	Course
Smith	353
Smith	379
Smith	221
Clark	353
Clark	351
Clark	379
Clark	456
Turner	353
Turner	456
Turner	272
Jamieson	353
Jamieson	379
Jamieson	221
Jamieson	456
Jamieson	469

COURSE	
Course	Dept
353	Comp Sci
379	Comp Sci
221	Decision Sci
351	Comp Sci
456	Mathematics
272	Chemistry
469	Mathematics

FACULTY	
Prof	Dept
Smith	Comp Sci
Clark	Comp Sci
Turner	Chemistry
Jamieson	Mathematics

2NF Drawback example:

TEACHES				
Course	Prof	Room	Room_Cap	Enrol_I_m
353	Smith	A532	45	40
351	Smith	C320	100	60
355	Clark	H940	400	300
456	Turner	B278	50	45
459	Jamieson	D110	50	45

Third Normal Form (3NF) (cont)

Definition (Third Normal Form): A relation scheme is in *third normal form (3NF)* if for all non-trivial FD's of the form $X \rightarrow A$ that hold, either X is a superkey or A is a prime attribute.

Note: a FD $X \rightarrow Y$ is trivial iff Y is a subset of X .

Alternative definition: A relation scheme is in third normal form if every non-prime attribute is fully functionally dependent on the keys and not transitively dependent on any key.

A database scheme is in 3NF if all its relations are in 3NF.

TEACHES can be decomposed into 3NF:

ROOM_DETAILS		
Room	Room_Cap	Enrol_I_mt
A532	45	40
C320	100	60
B278	50	45
D110	50	45
H940	400	300

COURSE_DETAILS		
Course	Prof	Room
353	Smith	A532
351	Smith	C320
456	Turner	B278
459	Jamieson	D110
355	Clark	H940

Another example:

LOTS

Property_Id	City	Lot_No	Area	Price	Tax_Rate

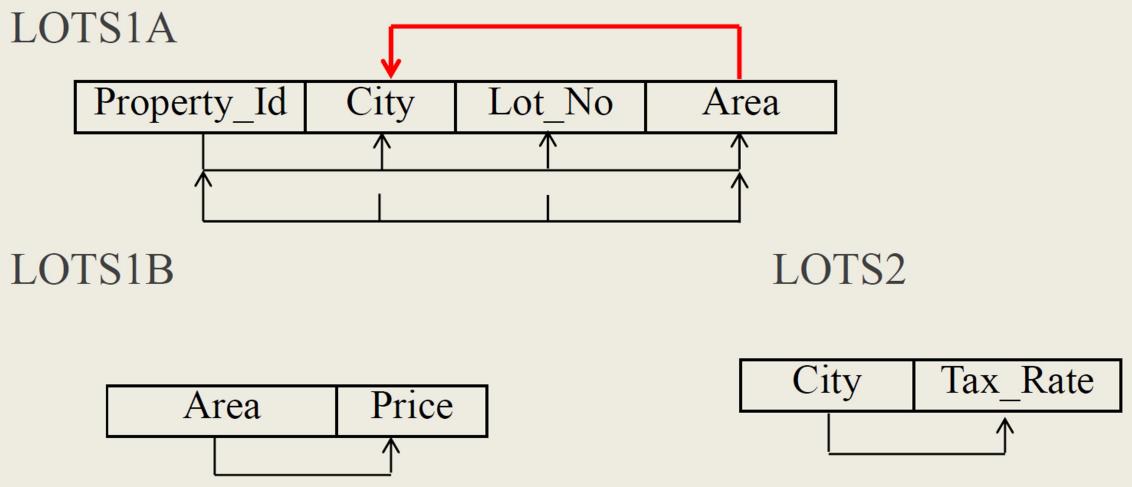
LOTS1

City	Tax_Rate

LOTS2

Property_Id	City	Lot_No	Area	Price

We could fix this too:



Boyce-Codd Normal Form (BCNF)

Definition (Boyce-Codd Normal Form):

A relation scheme is in *Boyce-Codd Normal Form* (BCNF) if whenever $X \rightarrow A$ holds and $X \rightarrow A$ is non-trivial, X is a superkey.

A database scheme is in BCNF if all its relations are in BCNF.

We can make our example into BCNF:

LOTS1AA

Property_Id	Lot_No	Area

LOTS1AB

Area	City

LOTS1B

Area	Price

LOTS2

City	Tax_Rate

LOTS1AA

Property_Id	Lot_No	Area

LOTS1AB

Area	City	Price

LOTS2

City	Tax_Rate

Relational Database Design

Anomalies can be removed from relation designs by decomposing them until they are in a normal form.

Several problems should be investigated regarding a decomposition.

A decomposition of a relation scheme, R, is a set of relation schemes $\{R_1, \dots, R_n\}$ such that $R_i \subseteq R$ for each i, and $\bigcup_{i=1}^n R_i = R$

Note that in a decomposition $\{R_1, \dots, R_n\}$, the intersect of each pair of R_i and R_j does not have to be empty.

Example: $R = \{A, B, C, D, E\}$, $R_1 = \{A, B\}$, $R_2 = \{A, C\}$, $R_3 = \{C, D, E\}$

A naive decomposition: each relation has only attribute.

A good decomposition should have the following two properties.

Dependency Preserving

Examples

$F = \{ A \rightarrow BC, D \rightarrow EG, M \rightarrow A \}$, $R = (A, B, C, D, E, G, M, A)$

1) Given $R_1 = (A, B, C, M)$ and $R_2 = (C, D, E, G)$,

$F_1 = \{ A \rightarrow BC, M \rightarrow A \}$, $F_2 = \{ D \rightarrow EG \}$

$F = F_1 \cup F_2$. thus, dependency preserving

2) Suppose that $F' = F \cup \{ M \rightarrow D \}$. R_1 and R_2 remain the same.

Thus, F_1 and F_2 remain the same.

We need to verify if $M \rightarrow D$ is inferred by $F_1 \cup F_2$.

Since $M^+ |_{F_1 \cup F_2} = \{M, A, B, C\}$, $M \rightarrow D$ is not inferred by $F_1 \cup F_2$.

Thus, R_1 and R_2 are not dependency preserving regarding F' .

3) $F'' = \{ A \rightarrow BC, D \rightarrow EG, M \rightarrow A, M \rightarrow C, C \rightarrow D, M \rightarrow D \}$

$F_1 = \{ A \rightarrow BC, M \rightarrow A, M \rightarrow C \}$, $F_2 = \{ D \rightarrow EG, C \rightarrow D \}$

It can be verified that $M \rightarrow D$ is inferred by F_1 and F_2 .

Thus, $F''^+ = (F_1 \cup F_2)^+$

Hence, R_1 and R_2 are dependency preserving regarding F'' .

Lossless Join Decomposition

A second necessary property for decomposition:

A decomposition $\{R_1, \dots, R_n\}$ of R is a *lossless join* decomposition with respect to a set F of FD's if for every relation instance r that satisfies F :

$$r = \pi R_1(r) \bowtie \dots \bowtie \pi R_n(r).$$

If $r \subset \pi R_1(r) \bowtie \dots \bowtie \pi R_n(r)$, the decomposition is *lossy*.

Example

STUDENT ADVISED BY

Name	Department	Advisor
Jones	Comp Sci	Smith
Ng	Chemistry	Turner
Martin	Physics	Bosky
Dulles	Decision Sci	Hall
Duke	Mathematics	James
James	Comp Sci	Clark
Evan	Comp Sci	Smith
Baxter	English	Bronte

STUDENT DEPARTMENT

Name	Department
Jones	Comp Sci
Ng	Chemistry
Martin	Physics
Duke	Mathematics
Dulles	Decision Sci
James	Comp Sci
Evan	Comp Sci
Baxter	English

DEPARTMENT ADVISED BY

Department	Advisor
Comp Sci	Smith
Chemistry	Turner
Physics	Bosky
Decision Sci	Hall
Mathematics	James
Comp Sci	Clark
English	Bronte

If we join these decomposed relations we get:

Name	Department	Advisor
Jones	Comp Sci	Smith
Jones	Comp Sci	Clark*
Ng	Chemistry	Turner
Martin	Physics	Bosky
Dulles	Decision Sci	Hall
Duke	Mathematics	James
James	Comp Sci	Smith*
James	Comp Sci	Clark
Evan	Comp Sci	Smith
Evan	Comp Sci	Clark*
Baxter	English	Bronte

Testing for the lossless join property

Example 4: $R = (A, B, C, D, E, G)$,

$$F = \{AB \rightarrow G, C \rightarrow DE, A \rightarrow B\}.$$

Let $R_1 = (A, B)$, $R_2 = (C, D, E)$ and $R_3 = (A, C, G)$.

A	B	C	D	E	G	A	B	C	D	E	G	A	B	C	D	E	G
R_1	a	a	b	b	b	R_1	a	a	b	b	b	R_1	a	a	b	b	b
R_2	b	b	a	a	a	R_2	b	b	a	a	a	R_2	b	b	a	a	b
R_3	a	b	a	b	b	R_3	a	b	a	x	x	R_3	a	x	a	a	a

Lossless decomposition into BCNF

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

$$R1 = (C, D, E, G), R2 = (A, B, C, D)$$

$$R11 = (C, E, G), R12 = (E, D) \text{ due to } E \rightarrow D$$

$$R21 = (A, B, C), R22 = (C, D) \text{ because of } C \rightarrow D$$

Example: (From Desai 6.31)

Find a BCNF decomposition of the relation scheme below:

SHIPPING(Ship , Capacity , Date , Cargo , Value)

F consists of:

$Ship \rightarrow Capacity$

$\{Ship, Date\} \rightarrow Cargo$

$\{Cargo, Capacity\} \rightarrow Value$

From $Ship \rightarrow Capacity$, we decompose *SHIPPING* into

$R_1(Ship, Date, Cargo, Value)$

Key: $\{Ship, Date\}$

A nontrivial FD in F^+ violates BCNF:

$\{Ship, Cargo\} \rightarrow Value$

and

$R_2(Ship, Capacity)$

Key: $\{Ship\}$

$Ship \rightarrow Capacity$ $\{Ship, Date\} \rightarrow Cargo$ $\{Cargo, Capacity\} \rightarrow Value$
--

Only one nontrivial FD in F^+ : $Ship \rightarrow Capacity$

R_1 is not in BCNF so we must decompose it further into

$R_{11}(Ship, Date, Cargo)$

Key: $\{Ship, Date\}$

Ship → Capacity
$\{Ship, Date\} \rightarrow Cargo$
$\{Cargo, Capacity\} \rightarrow Value$

Only one nontrivial FD in F^+ with single attribute on the right side: $\{Ship, Date\} \rightarrow Cargo$

and

$R_{12}(Ship, Cargo, Value)$

Key: $\{Ship, Cargo\}$

Only one nontrivial FD in F^+ with single attribute on the right side:
 $\{Ship, Cargo\} \rightarrow Value$

This is in BCNF and the decomposition is lossless but not dependency preserving (the FD $\{Capacity, Cargo\} \rightarrow Value$) has been lost.

Or we could have chosen $\{Cargo, Capacity\} \rightarrow Value$, which would give us:

$R_1(Ship, Capacity, Date, Cargo)$

Key: $\{Ship, Date\}$

A nontrivial FD in F^+ violates BCNF:

Ship → Capacity
$\{Ship, Date\} \rightarrow Cargo$
$\{Cargo, Capacity\} \rightarrow Value$

$Ship \rightarrow Capacity$

and

$R_2(Cargo, Capacity, Value)$

Key: $\{Cargo, Capacity\}$

Only one nontrivial FD in F^+ with single attribute on the right side: $\{Cargo, Capacity\} \rightarrow Value$

and then from $Ship \rightarrow Capacity$,

$R_{11}(Ship, Date, Cargo)$

Key: $\{Ship, Date\}$

$Ship \rightarrow Capacity$
$\{Ship, Date\} \rightarrow Cargo$
$\{Cargo, Capacity\} \rightarrow Value$

Only one nontrivial FD in F^+ with single attribute

on the right side: $\{Ship, Date\} \rightarrow Cargo$

And

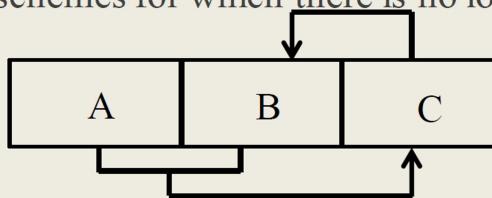
$R_{12}(Ship, Capacity)$

Key: $\{Ship\}$

Only one nontrivial FD in F^+ : $Ship \rightarrow Capacity$

This is in BCNF and the decomposition is both lossless and dependency preserving.

However, there are relation schemes for which there is no lossless, dependency preserving decomposition into BCNF.



Lossless and dependency-preserving decomposition into 3NF

A lossless and dependency-preserving decomposition into 3NF is always possible.

More definitions regarding FD's are needed.

A set F of FD's is minimal if

1. Every FD $X \rightarrow Y$ in F is simple: Y consists of a single attribute,
2. Every FD $X \rightarrow A$ in F is *left-reduced*: there is no proper subset $Y \subset X$ such that $X \rightarrow A$ can be replaced with $Y \rightarrow A$.

that is, there is no $Y \subset X$ such that

\rightarrow Iff $F \models Y \rightarrow A$

$$((F - \{X \rightarrow A\}) \cup \{Y \rightarrow A\})^+ = F^+$$

3. No FD in F can be removed; that is, there is no FD $X \rightarrow A$ in F

such that

\rightarrow Iff $X \rightarrow A$ is inferred
From $F - \{X \rightarrow A\}$

$$(F - \{X \rightarrow A\})^+ = F^+.$$

Example:

$$R = (A, B, C, D, E, G)$$

$$F = \{A \rightarrow BCD, B \rightarrow CDE, AC \rightarrow E\}$$

$$\text{Step 1: } F' = \{A \rightarrow B, A \rightarrow C, A \rightarrow D, B \rightarrow C, B \rightarrow D, B \rightarrow E, AC \rightarrow E\}$$

$$\text{Step 2: } AC \rightarrow E$$

$$C^+ = \{C\}; \text{ thus } C \rightarrow E \text{ is not inferred by } F'.$$

$$\text{Hence, } AC \rightarrow E \text{ cannot be replaced by } A \rightarrow E.$$

$$A^+ = \{A, B, C, D, E\}; \text{ thus, } A \rightarrow E \text{ is inferred by } F'.$$

$$\text{Hence, } AC \rightarrow E \text{ can be replaced by } A \rightarrow E.$$

$$F'' = \{A \rightarrow B, A \rightarrow C, A \rightarrow D, A \rightarrow E, B \rightarrow C, B \rightarrow D, B \rightarrow E\}$$

$$\text{Step 3: } A^+|_{F''} - \{A \rightarrow B\} = \{A, C, D, E\}; \text{ thus } A \rightarrow B \text{ is not inferred by } F'' - \{A \rightarrow B\}.$$

That is, $A \rightarrow B$ is not redundant.

$$A^+|_{F''} - \{A \rightarrow C\} = \{A, B, D, E\}; \text{ thus, } A \rightarrow C \text{ is redundant.}$$

Thus, we can remove $A \rightarrow C$ from F'' to obtain F''' .

Iteratively, we can $A \rightarrow D$ and $A \rightarrow E$ but not the others.

$$\text{Thus, } F_{\min} = \{A \rightarrow B, B \rightarrow C, B \rightarrow D, B \rightarrow E\}.$$

3NF decomposition algorithm

Algorithm 3NF decomposition

1. Find a minimum cover F' of F .

2. For each left side X that appears in F' , do:

create a relation schema $X \cup A_1 \cup A_2 \dots \cup A_m$ where $X \rightarrow \{A_1\}, \dots, X \rightarrow \{A_m\}$ are all the

dependencies in F' with X as left side.

3. if none of the relation schemas contains a key of R ,

create one more relation schema that contains attributes that form a key for R .

Example:

$$R = (A, B, C, D, E, G)$$

$$F_{\min} = \{A \rightarrow B, B \rightarrow C, B \rightarrow D, B \rightarrow E\}.$$

Candidate key: (A, G)

$$R_1 = (A, B), R_2 = (B, C, D, E)$$

$$R_3 = (A, G)$$

Example: (From Desai 6.31)

Beginning again with the *SHIPPING* relation. The functional dependencies already form a canonical cover.

- From $Ship \rightarrow Capacity$, derive $R_1(\underline{Ship}, Capacity)$,

- From $\{Ship, Date\} \rightarrow Cargo$, derive

$$R_2(\underline{Ship}, \underline{Date}, Cargo),$$

- From $\{Capacity, Cargo\} \rightarrow Value$, derive

$$R_3(\underline{Capacity}, \underline{Cargo}, Value).$$

- There are no attributes not yet included and the original key $\{Ship, Date\}$ is included in R_2 .

Another Example: Apply the algorithm to the LOTS example given earlier.

A minimal cover is

$$\begin{aligned} & \{ Property_Id \rightarrow Lot_No, \\ & \quad Property_Id \rightarrow Area, \{City, Lot_No\} \rightarrow Property_Id, \\ & \quad Area \rightarrow Price, Area \rightarrow City, City \rightarrow Tax_Rate \}. \end{aligned}$$

This gives the decomposition:

$$R_1(\underline{Property_Id}, Lot_No, Area)$$

$$R_2(\underline{City}, \underline{Lot_No}, \underline{Property_Id})$$

$$R_3(\underline{Area}, Price, City)$$

$$R_4(\underline{City}, Tax_Rate)$$