

Design Theory for Relational Databases

Normalization

Guillaume Raschia — Polytech Nantes; Université de Nantes

Last update: January 16, 2024

Contents

Bad Design

Boyce-Codd Normal Form

Third Normal Form

Minimal Cover and 3NF Synthesis

The Many Other Normal Forms

MVD's and 4NF

[Source : J. Ullman, Stanford]

Bad Design

Goal of relational schema design is to avoid **anomalies** and **redundancy**.

- **Update anomaly**: one occurrence of a fact is changed, but not all occurrences
- **Insertion anomaly**: related facts are required when a tuple is inserted
- **Deletion anomaly**: valid fact is lost when a tuple is deleted

Example of Bad Design

Drinkers(name, addr, beersLiked, brewery, favBeer)

name	addr	beersLiked	brewery	favBeer
Alice	Nantes	Trompe Souris	La Divatte	Titan
Alice	???	Titan	Bouffay	???
Bob	Rennes	Titan	???	Titan

Data is redundant, because each of the ???'s can be figured out by using the FD's
name → **addr** **favBeer** and **beersLiked** → **brewery**

Bad Design: Update Anomalies

name	addr	beersLiked	brewery	favBeer
Alice	Nantes Nantes Vannes	Trompe Souris	La Divatte Bouffay	Titan
Alice	Nantes	Titan	Bouffay	Titan
Bob	Rennes	Titan	Bouffay	Titan

- If Alice moves to Vannes, will we remember to change each of her tuples?

Bad Design: Deletion Anomalies

name	addr	beersLiked	brewery	favBeer
Alice	Nantes	Trompe Souris	La Divatte	Titan
Alice	Nantes	Titan	Bouffay	Titan
Bob	Rennes	Titan	Bouffay	Titan

- If nobody likes Trompe Souris anymore, we lose track of the fact that La Divatte brews Trompe Souris

Bad Design: Insertion Anomalies

name	addr	beersLiked	brewery	favBeer
Alice	Nantes	Trompe Souris	La Divatte	Titan
Alice	Nantes	Titan	Bouffay	Titan
Bob	Rennes	Titan	Bouffay	Titan
Charlie	Nantes	Mistral	Aerofab	Mistral

- If Charlie comes into play, one must know beers s/he likes and their breweries, otherwise `null` values

BCNF

Definition (BCNF)

We say a relation R is in BCNF if whenever $X \rightarrow Y$ is a nontrivial FD that holds in R , X is a superkey

- Remember: **nontrivial** means Y is not contained in X
- Remember: a **superkey** is any superset of a key (not necessarily a proper superset)

Example: BCNF

Drinkers(name, addr, beersLiked, brewery, favBeer)

- FD's:
 - name \rightarrow addr favBeer
 - beersLiked \rightarrow brewery
- Only key is {name, beersLiked}
- In each FD, the left-hand side is not a superkey
- Any one of these FD's shows **Drinkers** is not in BCNF

Another Example

Beers(name, brewery, brewAddr)

- FD's:
 - name \rightarrow brewery
 - brewery \rightarrow brewAddr
- Only key is {name}
- name \rightarrow brewery does not violate BCNF, but brewery \rightarrow brewAddr does

Decomposition into BCNF

Given: relation R with FD's \mathcal{F}

1. Look among the given FD's for a BCNF violation $X \rightarrow Y$
 - If any FD following from \mathcal{F} violates BCNF, then there will surely be an FD in \mathcal{F} itself that violates BCNF
2. Compute X^+
 - Not all attributes, or else X is a superkey

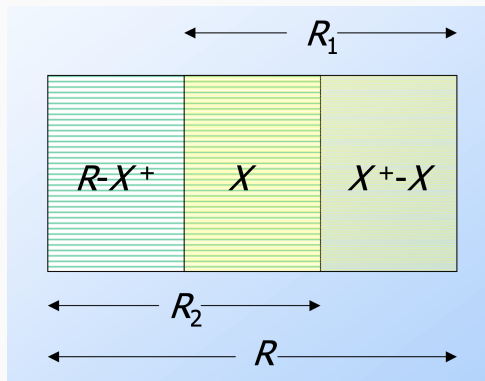
Decompose R Using $X \rightarrow Y$

3. Replace R by relations with schemas:

- $R_1 = X^+$
- $R_2 = R - (X^+ - X) = \overline{X^+} \cup X$

4. Project given FD's \mathcal{F} onto the two new relations

Decomposition Picture



J. Ullman

Example: BCNF Decomposition

Drinkers(name, addr, beersLiked, brewery, favBeer)

$\mathcal{F} = \{\text{name} \rightarrow \text{addr}, \text{name} \rightarrow \text{favBeer}, \text{beersLiked} \rightarrow \text{brewery}\}$

1. Pick BCNF violation $\text{name} \rightarrow \text{addr}$
2. Close the left-hand side: $\{\text{name}\}^+ = \{\text{name}, \text{addr}, \text{favBeer}\}$
3. Decomposed relations:
 - Drinkers1(name, addr, favBeer)
 - Drinkers2(name, beersLiked, brewery)

Example (cont'd)

We are not done; we need to check **Drinkers1** and **Drinkers2** for BCNF

- Projecting FD's is easy here
- For **Drinkers1**(name, addr, favBeer), relevant FD's are name \rightarrow addr and name \rightarrow favBeer
- Thus, {name} is the only key and **Drinkers1** is in BCNF

Example (cont'd)

- For `Drinkers2(name, beersLiked, brewery)`, the only FD is `beersLiked → brewery`, and the only key is {name, beersLiked}
- Violation of BCNF
- $\text{beersLiked}^+ = \{\text{beersLiked}, \text{brewery}\}$, so we decompose **Drinkers2** into:
 - `Drinkers3(beersLiked, brewery)`
 - `Drinkers4(name, beersLiked)`

The resulting decomposition of **Drinkers**

- **Drinkers1**(name, addr, favBeer)
- **Drinkers3**(beersLiked, brewery)
- **Drinkers4**(name, beersLiked)

Notice: **Drinkers1** tells us about drinkers, **Drinkers3** tells us about beers, and **Drinkers4** tells us the relationship between drinkers and the beers they like

3NF

Third Normal Form — Motivation

There is one structure of FD's that causes trouble when we decompose

- $AB \rightarrow C$ and $C \rightarrow B$
 - Example: A = street address, B = city, C = zip code
- There are two keys, $\{A,B\}$ and $\{A,C\}$
- $C \rightarrow B$ is a BCNF violation, so we must decompose into AC, BC

We Cannot Enforce FD's

The problem is that if we use AC and BC as our database schema, we cannot enforce the FD $AB \rightarrow C$ by checking FD's in these decomposed relations

A = street, B = city, and C = zip

$R =$	<table><tr><th>street</th><th>zip</th></tr><tr><td>50 Otages</td><td>44000</td></tr><tr><td>50 Otages</td><td>44100</td></tr></table>	street	zip	50 Otages	44000	50 Otages	44100
street	zip						
50 Otages	44000						
50 Otages	44100						

$S =$	<table><tr><th>zip</th><th>city</th></tr><tr><td>44000</td><td>Nantes</td></tr><tr><td>44100</td><td>Nantes</td></tr></table>	zip	city	44000	Nantes	44100	Nantes
zip	city						
44000	Nantes						
44100	Nantes						

{street, zip} is key in R , {zip} is key in S

An Unenforceable FD

$$R \bowtie S =$$

street	city	zip
50 Otages	Nantes	44000
50 Otages	Nantes	44100

$R \bowtie S$ joins tuples with equal zip codes

Although no FD's were violated in the decomposed relations, FD **street city \rightarrow zip** is violated by the database as a whole

- 3rd Normal Form (3NF) modifies the BCNF condition so we do not have to decompose in this problem situation

Definition (3NF (C. Zaniolo, 1982¹))

Every FD $X \rightarrow A$ satisfies one of those three conditions:

1. $X \rightarrow A$ is trivial
2. X is a superkey
3. A is prime (more flexible than BCNF)
 - An attribute is prime if it is a member of any key

In other words, a nontrivial FD $X \rightarrow A$ violates 3NF if and only if X is not a superkey, or A is not prime

¹Equivalent to (E.F. Codd, 1971).

Example: 3NF

- In our problem situation with FD's $AB \rightarrow C$ and $C \rightarrow B$, we have keys AB and AC
- A , B and C are each prime
- Although $C \rightarrow B$ violates BCNF, it **does not violate 3NF** (B is prime)
- One can decide not to decompose to BCNF, still being 3NF

What 3NF and BCNF Give You?

Two important properties of a decomposition

1. **Lossless Join:** it should be possible to project the original relations onto the decomposed schema, and then reconstruct the original
2. **Dependency Preservation:** it should be possible to check in the projected relations whether all the given FD's are satisfied

3NF and BCNF (cont'd)

- We can get (1) with a **BCNF** decomposition
- We can get both (1) and (2) with a **3NF** decomposition
- But we can't always get (1) and (2) with a **BCNF** decomposition
 - street-city-zip is an example

Testing for a Lossless Join

- If we project R onto R_1, R_2, \dots, R_n , can we recover R by rejoining?
 - A projected fragment: $R_i = \pi_{X_i}(R)$
 - Does R equal to $R_1 \bowtie R_2 \bowtie \dots \bowtie R_n$?
- Any tuple in R can be recovered from its projected fragments
 - $R \subseteq R_1 \bowtie R_2 \bowtie \dots \bowtie R_n$ is obvious
- So the only question is: when we rejoin, do we ever get back something we didn't have originally?
 - $R \supseteq R_1 \bowtie R_2 \bowtie \dots \bowtie R_n$ must be proved

The Chase Test

- Suppose tuple t comes back in the join
- Then t is the join of projections of some tuples of R , one for each R_i of the decomposition
- Can we use the given FD's to show that one of these tuples must be t ?

Procedure

1. Start by assuming $t = abc \dots$
2. For each i , there is a tuple s_i of R that has a, b, c, \dots in the attributes of R_i
3. s_i can have any values in other attributes
4. We'll use the same letter as in t , but with a subscript, for these components

Example: The Chase

- Let $R(A, B, C, D)$ and the decomposition be $R_1(A, B)$, $R_2(B, C)$ and $R_3(C, D)$
- Let the given FD's be $\mathcal{F} = \{C \rightarrow D, B \rightarrow A\}$
- Suppose the tuple $t = abcd$ is the join of tuples projected onto AB, BC, CD

The Tableau

Let's build an instance of R from tuple $t = abcd$ in $\pi_{AB}(R) \bowtie \pi_{BC}(R) \bowtie \pi_{CD}(R)$

R	A	B	C	D
from AB part of t	a	b	c_1	d_1
from BC part of t	a_2 a	b	c	d_2 d
from CD part of t	a_3	b_3	c	d

- Use $B \rightarrow A$ to state a_2 must be a
- Use $C \rightarrow D$ to state d_2 must be d

We've proved the second tuple must be $t = abcd$; then

$$(\pi_{AB}(R) \bowtie \pi_{BC}(R) \bowtie \pi_{CD}(R)) \subseteq R!$$

Build the Tableau, then

1. If two rows agree in the left side of a FD, make their right sides agree too
2. Always replace a subscripted symbol by the corresponding unsubscripted one, if possible
3. If we ever get an unsubscripted row, we know any tuple in the project-join is in the original table (the join is **lossless**)
4. Otherwise, the final tableau is a **counterexample**

Example: Lossy Join

- Same relation $R(A, B, C, D)$ and same decomposition AB, BC, CD
- But with only the FD $C \rightarrow D$

R	A	B	C	D
$t_1 =$	a	b	c_1	d_1
$t_2 =$	a_2	b	c	d_2 d
$t_3 =$	a_3	b_3	c	d

- Use $C \rightarrow D$ to state d_2 must be d , and that's all

These three tuples are an example of R that shows **the join is lossy**:

- $abcd$ is not in R , but we can project and rejoin to get $t = abcd$

$$t_1[AB] \bowtie t_2[BC] \bowtie t_3[CD] = abcd$$

3NF Decomposition

There is always a **lossless-join and dependency-preserving 3NF** decomposition

How to achieve 3NF?

1. Perform the iterative binary decomposition process up to 3NF only
2. Use the **3NF Synthesis**
 - Need a **minimal basis** for the FD's

Minimal Cover and 3NF Synthesis

also known as **minimal basis** or even **canonical cover**

A set of FD's is a minimal cover iff

1. *rhs's* are single attributes
2. No **redundant FD**, ie FD that can be discarded
3. No **extraneous attribute in lhs**, ie that can be removed from the *lhs*

Constructing a Minimal Cover

Given a set of FD's \mathcal{F}

Finding a minimal cover \mathcal{F}_{\min} requires:

1. Split *rhs*'s
2. Repeatedly try to remove an FD $X \rightarrow A$ and see if the remaining FD's are equivalent to the original
3. Repeatedly try to remove an attribute B from a *lhs* $BX \rightarrow A$ and see if the resulting FD's are equivalent to the original
4. Iterate 2-3 up to stable

Constructing a Minimal Cover — for Real

How to achieve Step 2: $\mathcal{F} - \{X \rightarrow A\} \equiv \mathcal{F}$?

$\Leftrightarrow \mathcal{F} - \{X \rightarrow A\} \models X \rightarrow A$? (to prove $\mathcal{F} - \{X \rightarrow A\}$ implies \mathcal{F})

\Leftrightarrow Check for $A \in X^+$ wrt $\mathcal{F} - \{X \rightarrow A\}$

How to achieve Step 3: $(\mathcal{F} - \{BX \rightarrow A\} \cup \{X \rightarrow A\}) \equiv \mathcal{F}$?

$\Leftrightarrow \mathcal{F}$ implies $(\mathcal{F} - \{BX \rightarrow A\} \cup \{X \rightarrow A\})$

• The other way round is obvious since $X \rightarrow A \models BX \rightarrow A$

$\Leftrightarrow \mathcal{F} \models X \rightarrow A$

\Leftrightarrow Check for $A \in X^+$ wrt \mathcal{F}

Example: Minimal Cover

Given the FD's $\mathcal{F} = \{ABC \rightarrow CD, A \rightarrow B, C \rightarrow A\}$

- Step 1: $\mathcal{F}_1 = \{ABC \rightarrow C, ABC \rightarrow D, A \rightarrow B, C \rightarrow A\}$
- Step2:
 - remove (trivial) $ABC \rightarrow C$ to get $\mathcal{F}_2 = \mathcal{F}_1 - \{ABC \rightarrow C\}$
 - cannot remove $ABC \rightarrow D$, since $D \notin ABC_{\mathcal{F}_2 - \{ABC \rightarrow D\}}^+ = ABC$
 - cannot remove $A \rightarrow B$, since $B \notin A_{\mathcal{F}_2 - \{A \rightarrow B\}}^+ = A$
 - cannot remove $C \rightarrow A$, since $A \notin C_{\mathcal{F}_2 - \{C \rightarrow A\}}^+ = C$

Example: Minimal Cover (cont'd)

Given the FD's $\mathcal{F} = \{ABC \rightarrow CD, A \rightarrow B, C \rightarrow A\}$

Step 1-2 yields to $\mathcal{F}_2 = \{ABC \rightarrow D, A \rightarrow B, C \rightarrow A\}$

- Step 3:
 - remove A in $ABC \rightarrow D$ since $D \in BC_{\mathcal{F}_2}^+ = BCAD$. Define \mathcal{F}_3
 - remove B in $BC \rightarrow D$ since $D \in C_{\mathcal{F}_3}^+ = CABD$. Define \mathcal{F}_4
 - cannot remove C in $C \rightarrow D$ (singleton in lhs)
 - nothing to do for $A \rightarrow B$ and $C \rightarrow A$ (singletons in lhs)

Finally, $\mathcal{F}_{\min} = \mathcal{F}_4 = \{C \rightarrow D, A \rightarrow B, C \rightarrow A\}$

Properties of Minimal Cover

Given a set of FD's \mathcal{F}

- FD's in \mathcal{F}_{\min} are **irreducibles**
- $\mathcal{F}_{\min} \equiv \mathcal{F}$ that is equivalent to $\mathcal{F}_{\min}^+ = \mathcal{F}^+$
- \mathcal{F}_{\min} is **not unique**
 - depends on the nondeterministic choices in Steps 2-3
- \mathcal{F}_{\min} is **required for the 3NF synthetis** algorithm

3NF Synthesis

1. Create one relation for each *lhs* of FD's in the minimal cover
 - Schema is the union of *lhs* and set of *rhs*'s
2. Discard relation $R(X)$ if $S(XY)$ exists
3. If no key is contained in an FD, then add one relation whose schema is some key

Example: 3NF Synthesis

Relation $R(A, B, C, D, E)$ with FD's $\mathcal{F} = \{AB \rightarrow C, AB \rightarrow D, C \rightarrow B, E \rightarrow B\}$

Assume \mathcal{F} is a **minimal cover**

Decomposition

1. ABCD, CB and EB
2. Then, remove CB since $CB \subseteq ABCD$
3. And add AE for a key

Example: 3NF Synthesis (cont'd)

Resulting decomposition is $R_1(\underline{AB}CD)$, $R_2(\underline{E}B)$ and $R_3(\underline{A}\underline{E})$

- R_1 R_2 and R_3 are 3NF
- R_1 is not BCNF by $C \rightarrow B$

3NF Synthesis

- **Preserves dependencies:** each FD from a minimal cover is contained in a relation, thus preserved
- **Lossless Join:** use the chase to show that the row for the relation that contains a key can be made all unsubscripted variables
- **3NF:** hard part, a property of minimal covers

The Many Other Normal Forms

First Normal Form

The very baseline of Relational Database Design

1NF

Relation has (a) a **key** and (b) **atomic** columns and (c) **no repeating groups** of columns

- Sets or tuples or tables are not allowed as attribute values
- (beersLiked₁, beersLiked₂, BeersLiked₃) is not allowed as a subset of columns

2NF

1NF and every non-prime attribute is **fully** functionally dependent on keys

- Remind: non-prime attributes are not key attributes
- FD $X \rightarrow A$ is full *iff* A doesn't depend on a proper subset of X
- 2NF is not so relevant in DB design

Example: 2NF

Drinkers(name, addr, beersLiked, brewery, favBeer)

FD's are name \rightarrow addr favBeer and beersLiked \rightarrow brewery

- **Drinkers** is 1NF: (a) key, (b) atomic values (c) non-repeating attributes
- **Drinkers** is not 2NF
 - name beersLiked \rightarrow addr is not full, since name \rightarrow addr holds
- **Drinkers** cannot be 3NF either

1NF < 2NF < 3NF < BCNF

- Always try to decompose up to BCNF
- When one cannot get dependency preserving BCNF, may decide to stop to 3NF
- Denormalize below 3NF only for good reason (performance) or data model shift²

²nested relations, document db, key-value stores, ...

Beyond Functional Dependencies

- Multi-Valued Dependencies: 4NF
- FD + Join Dependencies: ETNF (H. Darwen et al., 2012)
- JD: 5NF, 6NF
- Domain and Key constraints: DKNF
- ...

MVD's and 4NF

Class Book

title

set of authors

publisher

set of keywords

- Straightforward to model in any programming language
- Tricky in relational database!

Basic Proposal

Either we ignore the normalization...

Title	Author	Publisher	Keyword
FoD	S. Abiteboul	Addison-Wesley	Database
FoD	R. Hull	Addison-Wesley	Database
FoD	V. Vianu	Addison-Wesley	Database
FoD	S. Abiteboul	Addison-Wesley	Logic
FoD	R. Hull	Addison-Wesley	Logic
FoD	V. Vianu	Addison-Wesley	Logic
TCB	J.D. Ullman	Pearson	Database
⋮	⋮	⋮	⋮

- Key: (Title, Author, Keyword)
- Not in 2NF, given Title \rightarrow Publisher

Intermediate State

...Or we go to 3NF, BCNF

Title	Publisher
FoD	Addison-Wesley

Title	Author	Keyword
FoD	S. Abiteboul	Database
FoD	R. Hull	Database
FoD	V. Vianu	Database
FoD	S. Abiteboul	Logic
FoD	R. Hull	Logic
FoD	V. Vianu	Logic

- But we still ignore the multivalued dependencies...

MVD's are full constraints on relation³

Definition (Multi-Valued Dependency)

Let R be a relation of schema $\{X, Y, Z\}$; $X \twoheadrightarrow Y$ holds whenever (x, y, z) and (x, t, u) both belong to R , it implies that (x, y, u) and (x, t, z) should also be in R

Informally

$X \twoheadrightarrow Y$ holds if for any value of X , there exists a well-defined set of values of Y and a well-defined set of values of Z , independent one with the other.

³All the attributes are necessarily involved.

MVD's (cont'd)

$$A \twoheadrightarrow B \mid C$$

A	B	C
a_1	b_1	c_1
a_1	b_1	c_2
a_1	b_2	c_1
a_1	b_2	c_2
a_2	b_1	c_1
a_2	b_1	c_3

MVD's from Physics:

- Department {Building} {Employee {Telephone}}
- MVD's = {Dpt \twoheadrightarrow Bding | Emp, Tel ; Dpt, Emp \twoheadrightarrow Tel | Bding}

Properties of MVD's

Inference System [FoD 1994, Theorem 8.3.5 p. 172]

Let U a given set of attributes and X, Y, Z are subsets of U ;

- Complementation: if $X \twoheadrightarrow Y$, then $X \twoheadrightarrow (U - Y)$
- Reflexivity: if $Y \subseteq X$, then $X \twoheadrightarrow Y$ (trivial)
- Augmentation: if $X \twoheadrightarrow Y$, then $XZ \twoheadrightarrow YZ$
- Transitivity: if $X \twoheadrightarrow Y$ and $Y \twoheadrightarrow Z$ then $X \twoheadrightarrow (Z - Y)$
- Conversion: if $X \twoheadrightarrow Y$, then $X \rightarrow Y$
- Interaction: if $X \twoheadrightarrow Y$ and $XY \rightarrow Z$, then $X \rightarrow (Z - Y)$

Armstrong's Axioms (both on FD's and MVD's) constitute a sound and complete inference system for the FD+MVD closure computation

Definition (4NF)

For every non trivial MVD $X \twoheadrightarrow Y$ in R , then X is a superkey

Straightforward extension of BCNF to MVD's.

Lossless-join decomposition of $R(X, Y, Z)$

Decomposition (X, Y) and (X, Z) is **lossless-join** iff $X \twoheadrightarrow Y$ holds in R

Back to the Class Book Running Example

Title	Publisher
FoD	Addison-Wesley

Title	Author	Keyword
FoD	S. Abiteboul	Database
FoD	R. Hull	Database
FoD	V. Vianu	Database
FoD	S. Abiteboul	Logic
FoD	R. Hull	Logic
FoD	V. Vianu	Logic

List of MVD's:

- Title \twoheadrightarrow Author | Keyword

The Ultimate Schema

...Up to the 4NF

Title	Publisher
FoD	Addison-Wesley

Title	Author
FoD	S. Abiteboul
FoD	R. Hull
FoD	V. Vianu

Title	Keyword
FoD	Database
FoD	Logic

Pros & Cons

- 3NF/BCNF design
 - must have
 - best trade-off between redundancy vs. decomposition
 - a priori BCNF, except on dependency loss, then design choice
- 4NF design
 - requires many joins in queries (performance pitfall)
 - and loses the big picture of class book entities
- 1NF relational view
 - eliminates the need for users/apps to perform deadly joins
 - but loses the one-to-one mapping between tuples and objects
 - has a large amount of redundancy
 - and could yield to insertion, deletion, update anomalies