Normalization

CS 4750 Database Systems

[A. Silberschatz, H. F. Korth, S. Sudarshan, Database System Concepts, Ch.7] [Ricardo and Urban, Database Illuminated, Ch.6] [https://www.w3schools.in/dbms/database-normalization/]

© Praphamontripong

Recap: FD

Consider a student_lecture relation

S_id	Address	Course	Teaching_assistant
1234	57 Hockanum Blvd	Database Systems	Minnie
2345	1400 E. Bellows	Database Systems	Humpty
3456	900 S. Detroit	Cloud Computing	Dumpty
1234	57 Hockanum Blvd	Web Programming Lang.	Mickey
5678	2131 Forest Lake Ln.	Software Analysis	Minnie

Assume that there is exactly one teaching assistant assigned to each student for every course

- 1. Determine all functional dependencies of the relation
- 2. Give an example of a super key and a candidate key

Recap: FD

Consider a student_lecture relation

S_id	Address	Course	Teaching_assistant
1234	57 Hockanum Blvd	Database Systems	Minnie
2345	1400 E. Bellows	Database Systems	Humpty
3456	900 S. Detroit	Cloud Computing	Dumpty
1234	57 Hockanum Blvd	Web Programming Lang.	Mickey
5678	2131 Forest Lake Ln.	Software Analysis	Minnie

1. Determine all functional dependencies of the relation

S_id → Address
S_id, Course → Teaching_assistant



Assume: there is exactly one teaching assistant assigned to each student for every course

Recap: FD

Consider a student_lecture relation

S_id	Address	Course	Teaching_assistant
1234	57 Hockanum Blvd	Database Systems	Minnie
2345	1400 E. Bellows	Database Systems	Humpty
3456	900 S. Detroit	Cloud Computing	Dumpty
1234	57 Hockanum Blvd	Web Programming Lang.	Mickey
5678	2131 Forest Lake Ln.	Software Analysis	Minnie

2. Give an example of a super key and a candidate key

The set of all the attributes is a trivial super key

(S_id, Course) is a candidate key

Assume: there is exactly one teaching assistant assigned to each student for every course

Recap: Attribute Closure

Given a relation R (A, B, C, D, E, F, G) with the following FDs FDs = { A \rightarrow BC, E \rightarrow CF, B \rightarrow E, CD \rightarrow EF, A \rightarrow G } Find the electron of A (A +)

Find the closure of A (A+)

$A \rightarrow B$ and $A \rightarrow C$	(decompose A \rightarrow BC)
$A \rightarrow E$	(transitive A \rightarrow B and B \rightarrow E)
$A \rightarrow CF$	(transitive A \rightarrow E and E \rightarrow CF)
$A \rightarrow G$	(FD given)
Thus, $A + = \{ ABCE \}$	=G }

General Design Guidelines

- Semantics of attributes should self-evident
- Avoid redundancy between tuples, relations
- Avoid NULL values in tuples
- If certain tuples should not exist, don't allow them

Database design = process or organizing data into a database model by considering data needs to be stored and the interrelationship of the data

> Database design is about characterizing data and the organizing data

How to describe properties we know or see in the data How to organize data to promote ease of use and efficiency

revisit

Normalization

- Normalization = technique of organizing data in a database
- Two purposes:
 - Eliminating redundant data
 - Avoid storing the same data in multiple tables
 - Ensuring data dependencies make sense
 - Store data logically only related data in a table, nothing else
- Need to refine schema

Schema Refinement

- Constraints, in particular functional dependencies, cause problems
- Must understand when and how constraints cause redundancy
- Refinement is needed when redundancy exists
- **Decomposition** main refinement technique
 - Example: replace ABCD with [AB and BCD] or [ACD and ABD]
 - Judgment call:
 - Is there a reason to decompose a relation?
 - What problems (if any) does the decomposition cause?

Decomposition

Suppose a relation R contains attribute A_1 , ..., A_n . A decomposition of R consists of replacing R by two or more relations such that

- Each new relation schema contains a subset of the attributes of *R* (and no attribute that do not appear in *R*)
- Every attribute of R appears as an attribute of at least one of the new relations

Three potential problems:

Tradeoff: must consider these issues vs. redundancy

- Some queries become more expensive
- Given instances of the decomposed relations, we may not be able to reconstruct the original relation
- Checking some dependencies may require joining the the decomposed relations

Properties of Decomposition

Lossless join

- Employee = R1 \bowtie R2 (\bowtie "natural join")
- No gain or loose columns / rows
- R1 \cap R2 \neq { }
- R1 \cap R2 \rightarrow R1 or R1 \cap R2 \rightarrow R2 (R1 \cap R2 is a super key of R1 or R2)

Dependency preserving

• Every dependency is in the same relation (thus, when checking a dependency, no need to join tables)

No redundancy

• For every nontrivial FD, a determinant must be a superkey (solved through the normal forms)

Lossless-Join Decomposition

computingID	name	year	hourly_rate	hours_worked
ht1y	Humpty	4	12	20
dt2y	Dumpty	3	10	20
md3y	Mickey	4	12	15
mn4e	Minnie	4	12	16
dh5h	Duhhuh	3	10	10

Employee = R1 \bowtie R2

No gain or loose columns

R1			
computingID	name	year	hours_worked
ht1y	Humpty	4	20
dt2y	Dumpty	3	20
md3y	Mickey	4	15
mn4e	Minnie	4	16
dh5h	Duhhuh	3	10

$$R1 \cap R2 \neq \{ \}$$
Super key
$$R2$$

$$R2$$

$$4$$

$$12$$

$$3$$

$$10$$

D

Lossless-Join Decomposition

R1						
computingID	name	year	hours_worked			
ht1y	Humpty	4	20			
dt2y	Dumpty	3	20			R2
md3y	Mickey	4	15		year	hourly_rate
mn4e	Minnie	4	16	\mathbf{M}	4	12
dh5h	Duhhuh	3	10	\mathbf{N}	3	10

Reconstruct the original relation No gain or loose columns / rows

	Employee				
	computingID	name	year	hourly_rate	hours_worked
	ht1y	Humpty	4	12	20
=	dt2y	Dumpty	3	10	20
	md3y	Mickey	4	12	15
	mn4e	Minnie	4	12	16
	dh5h	Duhhuh	3	10	10

Let's Try (1) Lossless Join decomposition?



Νο

Let's Try (2) Lossless Join decomposition?





Normalization

- **1NF**: each column is atomic, flat
- **2NF**: 1NF + no partial dependency -- [outdated]
- **3NF**: 2NF + lossless-join + dependency preserving -- [our focus]
- **BCNF**: 1NF + lossless join + redundancy free -- [our focus]
- **4NF:** no multi-valued dependency -- [out of CS 4750 scope]
- **5NF**: 4NF + cannot be further non loss decomposed -- [out of CS 4750 scope too complicated]
- **6NF**: 5NF + every join dependency is trivial -- [out of CS 4750 scope somewhat unrealistic]

First Normal Form (1NF)

- Every attribute/column has a single (atomic) value
- Values stored in a column should be of the same domain
- The order in which data is stored does not matter

com	putingID	na	name p		phone		department		
ht1y		Hum	pty	111-111-1111		(Computer Science, Math		
dt2y		Dumpty		222-222-2222		pty 222-222-2222			Biology
md3y	computin	gID	nam	ne	phone		department		
mn4e	¹ e ht1y Hun		Humpt	у	111-111-1111		Computer Science		
	ht1y		Humpt	У	111-111-1111		Math		
	dt2y		Dumpt	У	222-222-2222	2	Biology		
	md3y Mi		Mickey		333-333-3333	3	Computer Science		
	mn4e		Minnie		444-444-4444	ŀ	Computer Science		

Suppose we know that a student may be in multiple departments

Second Normal Form (2NF)

1NF + no partial dependency (FDs)



Suppose we also know course \rightarrow instructor

Let's simplify the name: R (A, B, C, D), AB is a candidate key FDs { AB \rightarrow C, AB \rightarrow D, B \rightarrow D }

Since B is part of a candidate key, D depends on a part of a key "Partial dependency"

Note: many-to-many relationship

Second Normal Form (2NF)

To convert the table into 2NF, decompose the table to remove partial dependency

	computingID		course		grade		instructor			
	ht	1y	cs1		B+		someone1			
	dt	2у	cs1		cs1		A-		someone1	
	dt	2у	cs2		А	someone2				
comp	utingID	course	grade	e A			someone1			
ht	:1y	cs1	B+		B	urse	instructor			
dt	:2y	cs1	A-		Δ	jui se	mstructor			
dt	:2v	cs2	А		· · · · ·	cs1	someone1			
m	d3y	cs1	А			cs2	someone2			
mi	n4e	cs2	В							
ma	d3y	cs2	A							

Non-key attributes must depend upon the whole of the candidate key

Third Normal Form (3NF)

- 2NF + lossless-join + dependency preserving
- For every non-trivial dependency, either LHS is a superkey or the RHS consists of prime attributes only
 No transitive

computingID	course	grade	textbook_id	title	
ht1y	cs1	B+	book1	Intro to Python	
dt2y	cs1	A-	book1	Intro to Python	
dt2y	cs2	А	book2	Intro to Java	$B \rightarrow D$?
md3y	cs1	А	book1	Intro to Python	
mn4e	cs2	В	book2	Intro to Java	
md3y	cs2	А	book2	Intro to Java)

Suppose we want to keep track of textbook_id and title for the course

Let's simplify the name: R (A, B, C, D, E), AB is a candidate key FDs { AB \rightarrow C, AB \rightarrow D, D \rightarrow E }

Since AB \rightarrow D, D \rightarrow E, D and E are non keys -- "transitive dependency" Problem with dependency ... Fix it!

dependencies

Third Normal Form (3NF)

To convert the table into 3NF, decompose the table to remove transitive dependency

	computi	ngID		course		grade	te	xtbo	ok_id		title	
	ht1y			cs1	B+		book1		Intro to Python		ython	
	dt2y			cs1	A-		book1		Intro to Python		ython	
	dt2v			cs2	А		book2		Intro to Java		Java	
com	nputingID	course	e	grade		course		tex	ktbook_id		, thon	
	ht1y	cs1		B+	cs1		s1		book1			ava
	dt2y	cs1		A-	CS		s2	b		book2		
	dt2y	cs2		А								ava
1	md3y	cs1		А			textbo		extbook_	ook_id		title
I	mn4e	cs2		В					book1		Intro	o to Pythor
	md3y	cs2	J	А					book2		Int	ro to Java

Ensure data integrity; no transitive dependency; dependency is in the same relation

Preserve all FDs but allow anomalies (may have redundancy)

Boyce-Codd Normal Form (BCNF)

- 1NF + lossless-join + redundant free
- For every non-trivial dependency, $X \rightarrow A$, X is a superkey.

	ţ	All dependencies must be from					
computingID	course	instructor	full key				
ht1y	cs1	someone1					
dt2y	cs1	someone1	Suppose we know that				
dt2y	cs2	someone2	instructor \rightarrow course				
md3y	cs1	someone3					
mn4e	cs2	someone2	Cannot have a non-key				
md3y	cs2	someone2	implies a key				

Let's simplify the name: R (A, B, C), AB is a candidate key FDs { AB \rightarrow C, C \rightarrow B }

Since C \rightarrow B, non-key implies a (part of) key

Boyce-Codd Normal Form (BCNF)

To convert the table into BCNF, decompose the table to remove non-keys that imply a key – to make all dependencies from a key



Remove redundant data; ensure data integrity; may have dependency across relations (need to join to check dependency) No transitive FDs, no non-key dependencies, but can lose FDs

Wrap-Up

- Properties of Decomposition
 - Lossless join
 - Dependency preserving
 - No redundancy
- Overview of normal forms

What's next?

- 3NF and decomposition
- BCNF and decomposition