

## 6.5830 Lecture 4

9.19.22

PS1 Due Next Time

Lab 1 due a week from Wednesday

### Relational Model Continued, and Schema Design and Normalization

Now that we understand relational algebra, we are almost ready to talk about how we actually implement a system that executes this algebra. But first, let's focus on how we design a database and choose a set of tables.

(break)

### **Schema Normalization**

Goal: Produce a collection of tables (schema) that is redundancy free

Q1 : Why is redundancy a problem?

- Because it leads to various anomalies when you try to update a table.
- Because it wastes space

<u>ss#</u>	<u>name</u>	<u>address</u>	<u>hobby</u>	<u>cost</u>
123	john	main st	dolls	\$
123	john	main st	bugs	\$
234	mary	lake st	bugs	\$
345	mary	lake st	tennis	\$\$
456	joe	first st	dolls	\$

What is the primary key? SS# + Hobby?

people have names and addresses, hobbies have costs

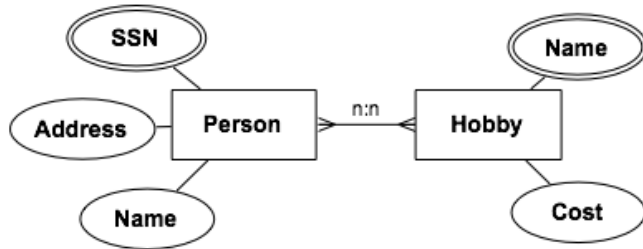
people can have multiple hobbies, and hobbies can be practiced by multiple people

Types of anomalies (Codd calls "inconsistencies")

- **Update anomaly** - change one address -- need to change the other
- **Insertion anomaly** - what if we want to add someone with no hobby?  
have to use a null?  
problem -- hobby is a part of the key!

Q2: What can we do to solve it? Normalize!

Most common way to do this is with an "entity relationship diagram"



## Entity Relationship Diagram

n:n relationships imply a mapping table from key of left table to key of right table

so we'd have

person (ssn, address, name)  
 hobby (name, cost)  
 personhobby (hobbyname, ssn)

1:n relationships imply a key-foreign key reference,

e.g., if every person has one hobby, we could add a hobbyid field

1:1 relationships can be merged into the same table (imply a strict dependencies) --

e.g., each person has a unique hobby, then their hobby can just be stored with them

(Decomposes into

ss#    name    address

1	joe	main st
2	jenny	lake st
3	jimmy	south st

hobby    cost

dolls	\$
bugs	\$
tennis	\$\$

ss#    hid

1	dolls
1	bugs
2	tennis
2	bugs
3	dolls)

we've eliminated the anomalies (only need to change one address, don't need nulls for hobbies, etc.)

what about:

ss# name address cost

hobby ss#

No redundancy, but we have lost some information. **"lossy decomposition"**

Let's formalize this idea a bit more to see why ER modeling leads to a good decomposition.

How do we do this systematically?

Let's understand where the redundancy comes from.

Looking at our hobbies example, SS# is sufficient to uniquely determine name and address, but the key of the table includes hobby, which means a given SS# can repeat, and hence so can the names and addresses for the SS#.

One way to think about this is in terms of "functions" -- in the sense that a given a particular input value, a function always produces the same output. So, e.g., a given SS# always produces the same name.

We write these like this:

SS# ---> Name

These kinds of relationships are called functional dependencies. Redundancy arises when the LHS of one the functional dependencies over the attributes in a table is not the key of the table.

For our tables above, we can write down some FDs:

FD1: SSN, Hobby -> Name, Address, Cost

FD2: SSN -> Name, Address

FD3: Hobby -> Cost

FD2 and FD3 sufficient to imply FD1 ==> "Armstrong's axioms"

Because for our "wide table", SSN is not a key, we can see that some info will be repeated each time an SSN is repeated with a hobby.

## Where do FDs come from?

Domain knowledge of database designer (not derived from data, though can check that data satisfies them!)

## Normal Forms

A table that has no redundancy is said to be in BCNF "Boyce Codd Normal Form"

Formally, a set of relations is in BCNF if:

For every functional dependency  $X \rightarrow Y$  in a set of functional dependencies  $F$  over relation  $R$   
 $X$  is a *superkey* key of  $R$ ,  
(where superkey means that  $X$  contains a key of  $R$ )

go to our hobbies example  
if we use the original example,

SSN  $\rightarrow$  Name, Address is not a superkey! So this is not in BCNF.

--> Redundancy

In non-decomposed hobbies schema Name, Addr repeated for each appearance of a given SSN

BCNF implies there is no redundant information -- e.g., that the association implied by any functional dependency is stored only once;

Observe that our schema after ER modeling is in BCNF (FDs for each table only have superkeys on the left side)

Decomposing into FD is easy -- just look at each FD, one by one, and check the conditions over each relation. If they don't apply to some relation  $R$ , split  $R$  into two relations,  $R_1$  and  $R_2$ , where  $R_1 = (X \cup Y)$  and  $R_2 = R - (X \cup Y)$ ,

Start with one "universal relation"

While some relation  $R$  is not in BCNF

Find an FD  $F=X \rightarrow Y$  that violates BCNF on  $R$

Split  $R$  into  $R_1 = (X \cup Y)$ ,  $R_2 = R - Y$

Example:

FD2: SSN  $\rightarrow$  Name, Address

FD3: Hobby  $\rightarrow$  Cost

R = S,N,A,H,C

R is not in BCNF, b/c of FD2 (N, A is not a primary key of R)

R1 = S,N,A, FD2 ; R2 = S,H,C, FD1', FD3

R2 not in BCNF, b/c of FD3

R3 = H, C (FD3)

R4 = S, H (FD1'')

Is it always possible to remove all redundancy?

Consider:

account, client, office:

Client, Office -> Account

Account -> Office

Account	Client	Office
a	joe	1
b	mary	1
a	john	1
c	joe	2

(key client,office)

Redundancy b/c the fact that account a is in office 1 is represented twice.

Not in BCNF b/c account is not a superkey of the table

This table is in "third normal form" (3NF) but not BCNF.

To put it into BCNF, we would have to decompose into

Account, Office

Client

But we've "lost" the first dependency now, in the sense that it can't be checked by looking at a single table.

3NF preserves all dependencies, but may have redundancy, and BCNF removes all redundancy but may drop some dependencies.

Comparable algorithms for 3NF, not going into details.

## **Denormalization**

When and where do we want to do it?

Do we always want to decompose a relation? Why or why not?

Generally speaking, decomposition :

- decreases storage overhead by eliminating redundancy and
- increases query costs by adding joins.

This isn't always true! Sometimes it increases storage overhead or decreases query costs.

Sometimes (for performance issues) you don't want to decompose.

### **So how much does this really matter?**

Eliminating redundancy really is important.  
Adding lots of joins can really screw performance.

These two are sometimes at odds with each other.

In practice, what people do is what we did for hobbies -- think about entities, join them on keys. "Entity relationship" model provides a way to do this and will result in something in BCNF.