

CompSci 516 Database Systems

Lecture 7 Relational Calculus (revisit) And Normal Forms

Instructor: Sudeepa Roy

Duke CS, Fall 2019

CompSci 516: Database Systems

1

Announcements

- **HW1 Deadlines!**
 - Today: parser and Q1-Q3
 - Q4: next Tuesday
 - Q5 (3 RA questions will be posted today): next Thursday
- **2 late days with penalty apply for individual deadlines**
 - If you are still parsing XML
 - Remember to start early next time from first day
 - HW2 and HW3 typically take more time and effort!

Duke CS, Fall 2019

CompSci 516: Database Systems

2

Today's topic

- Revisit RC
- Finish Normalization
- From Thursday: Database Internals

Acknowledgement:
The following slides have been created adapting the instructor material of the [RG] book provided by the authors Dr. Ramakrishnan and Dr. Gehrke, and with the help of slides by Dr. Magda Balazinska and Dr. Dan Suciu
CompSci 516: Database Systems

Duke CS, Fall 2019

3

Relational Calculus (RC) (Revisit from Lecture 4)

Duke CS, Fall 2019

CompSci 516: Database Systems

4

Logic Notations

- \exists There exists
- \forall For all
- \wedge Logical AND
- \vee Logical OR
- \neg NOT
- \Rightarrow Implies

TRC: example

Sailors(sid, sname, rating, age)
Boats(bid, bname, color)
Reserves(sid, bid, day)

- Find the name and age of all sailors with a rating above 7

\exists There exists

$\{P \mid \exists S \in \text{Sailors} (S.\text{rating} > 7 \wedge P.\text{sname} = S.\text{sname} \wedge P.\text{age} = S.\text{age})\}$

- P is a tuple variable
 - with exactly two fields sname and age (schema of the output relation)
 - $P.\text{sname} = S.\text{sname} \wedge P.\text{age} = S.\text{age}$ gives values to the fields of an answer tuple
- Use parentheses, \forall \exists \vee \wedge $>$ $<$ $=$ \neq \neg etc as necessary
- $A \Rightarrow B$ is very useful too

– next slide

Duke CS, Fall 2019

CompSci 516: Database Systems

6

$A \Rightarrow B$

- A “implies” B
- Equivalently, if A is true, B must be true
- Equivalently, $\neg A \vee B$, i.e.
 - either A is false (then B can be anything)
 - otherwise (i.e. A is true) B must be true

Duke CS, Fall 2019

CompSci 516: Database Systems

7

Useful Logical Equivalences

$$\bullet \forall x P(x) = \neg \exists x [\neg P(x)]$$

\exists	There exists
\forall	For all
\wedge	Logical AND
\vee	Logical OR
\neg	NOT

$$\bullet \neg(P \vee Q) = \neg P \wedge \neg Q$$

$$\bullet \neg(P \wedge Q) = \neg P \vee \neg Q$$

} de Morgan's laws

– Similarly, $\neg(\neg P \vee Q) = P \wedge \neg Q$ etc.

$$\bullet A \Rightarrow B = \neg A \vee B$$

Duke CS, Fall 2019

CompSci 516: Database Systems

8

TRC: example

Sailors(sid, sname, rating, age)
 Boats(bid, bname, color)
 Reserves(sid, bid, day)

- Find the names of sailors who have reserved at least two boats

Duke CS, Fall 2019

CompSci 516: Database Systems

9

TRC: example

Sailors(sid, sname, rating, age)
 Boats(bid, bname, color)
 Reserves(sid, bid, day)

- Find the names of sailors who have reserved at least two boats

$$\{P \mid \exists S \in \text{Sailors} (\exists R1 \in \text{Reserves} \exists R2 \in \text{Reserves} (S.\text{sid} = R1.\text{sid} \wedge S.\text{sid} = R2.\text{sid} \wedge R1.\text{bid} \neq R2.\text{bid}) \wedge P.\text{sname} = S.\text{sname})\}$$

Duke CS, Fall 2019

CompSci 516: Database Systems

10

TRC: example

Sailors(sid, sname, rating, age)
 Boats(bid, bname, color)
 Reserves(sid, bid, day)

- Find the names of sailors who have reserved all boats

Duke CS, Fall 2019

CompSci 516: Database Systems

11

TRC: example

Sailors(sid, sname, rating, age)
 Boats(bid, bname, color)
 Reserves(sid, bid, day)

- Find the names of sailors who have reserved all boats

$$\{P \mid \exists S \in \text{Sailors} [\forall B \in \text{Boats} (\exists R \in \text{Reserves} (S.\text{sid} = R.\text{sid} \wedge R.\text{bid} = B.\text{bid})) \wedge (P.\text{sname} = S.\text{sname})]\}$$

Duke CS, Fall 2019

CompSci 516: Database Systems

12

TRC: example

Sailors(sid, sname, rating, age)
 Boats(bid, bname, color)
 Reserves(sid, bid, day)

- Find the names of sailors who have reserved all red boats

How will you change the previous TRC expression?

Duke CS, Fall 2019 CompSci 516: Database Systems 13

TRC: example

Sailors(sid, sname, rating, age)
 Boats(bid, bname, color)
 Reserves(sid, bid, day)

- Find the names of sailors who have reserved all red boats
 $\{P \mid \exists S \in \text{Sailors} (\forall B \in \text{Boats} (B.\text{color} = \text{'red'} \Rightarrow (\exists R \in \text{Reserves} (S.\text{sid} = R.\text{sid} \wedge R.\text{bid} = B.\text{bid}))) \wedge P.\text{sname} = S.\text{sname})\}$

Recall that $A \Rightarrow B$ is logically equivalent to $\neg A \vee B$
 so \Rightarrow can be avoided, but it is cleaner and more intuitive

Duke CS, Fall 2019 CompSci 516: Database Systems 14

More Examples: RC

- The famous “Drinker-Beer-Bar” example!

UNDERSTAND THE DIFFERENCE IN ANSWERS
FOR ALL FOUR DRINKERS

Duke CS, Fall 2019 CompSci 516: Database Systems 15

Likes(drinker, beer)
 Frequent(drinker, bar)
 Serves(bar, beer)

Drinker Category 1

Find drinkers that frequent some bar that serves some beer they like.

...

Duke CS, Fall 2019 CompSci 516: Database Systems 16

Likes(drinker, beer)
 Frequent(drinker, bar)
 Serves(bar, beer)

Drinker Category 1

Find drinkers that frequent some bar that serves some beer they like.

$\{x \mid \exists F \in \text{Frequent} (F.\text{rinker} = x.\text{rinker} \wedge \exists S \in \text{Serves} \exists L \in \text{Likes} (F.\text{rinker} = L.\text{rinker}) \wedge (F.\text{bar} = S.\text{bar}) \wedge (S.\text{beer} = L.\text{beer}))\}$

Duke CS, Fall 2019 CompSci 516: Database Systems 17

Likes(drinker, beer)
 Frequent(drinker, bar)
 Serves(bar, beer)

Drinker Category 2

Find drinkers that frequent some bar that serves some beer they like.

$\{x \mid \exists F \in \text{Frequent} (F.\text{rinker} = x.\text{rinker} \wedge \exists S \in \text{Serves} \exists L \in \text{Likes} (F.\text{rinker} = L.\text{rinker}) \wedge (F.\text{bar} = S.\text{bar}) \wedge (S.\text{beer} = L.\text{beer}))\}$

Find drinkers that frequent only bars that serves some beer they like.

...

Free HW question hint!

Duke CS, Fall 2019 CompSci 516: Database Systems 18

Likes(drinker, beer)
Frequents(drinker, bar)
Serves(bar, beer)

Drinker Category 2

Find drinkers that frequent some bar that serves some beer they like.

$$\{x \mid \exists F \in \text{Frequents} (F.\text{drinker} = x.\text{drinker} \wedge \exists S \in \text{Serves} \exists L \in \text{Likes} (F.\text{drinker} = L.\text{drinker} \wedge (F.\text{bar} = S.\text{bar}) \wedge (S.\text{beer} = L.\text{beer})))\}$$

Find drinkers that frequent only bars that serve some beer they like.

$$\{x \mid \exists F \in \text{Frequents} (F.\text{drinker} = x.\text{drinker}) \wedge [\forall F1 \in \text{Frequents} (F.\text{drinker} = F1.\text{drinker}) \Rightarrow \exists S \in \text{Serves} \exists L \in \text{Likes} ((F1.\text{bar} = S.\text{bar}) \wedge (F1.\text{drinker} = L.\text{drinker}) \wedge (S.\text{beer} = L.\text{beer}))]\}$$

19

Likes(drinker, beer)
Frequents(drinker, bar)
Serves(bar, beer)

Drinker Category 3

Find drinkers that frequent some bar that serves some beer they like.

$$\{x \mid \exists F \in \text{Frequents} (F.\text{drinker} = x.\text{drinker} \wedge \exists S \in \text{Serves} \exists L \in \text{Likes} (F.\text{drinker} = L.\text{drinker}) \wedge (F.\text{bar} = S.\text{bar}) \wedge (S.\text{beer} = L.\text{beer})))\}$$

Find drinkers that frequent only bars that serve some beer they like.

$$\{x \mid \exists F \in \text{Frequents} (F.\text{drinker} = x.\text{drinker}) \wedge [\forall F1 \in \text{Frequents} (F.\text{drinker} = F1.\text{drinker}) \Rightarrow \exists S \in \text{Serves} \exists L \in \text{Likes} ((F1.\text{bar} = S.\text{bar}) \wedge (F1.\text{drinker} = L.\text{drinker}) \wedge (S.\text{beer} = L.\text{beer}))]\}$$

Find drinkers that frequent some bar that serves only beers they like.

...

20

Likes(drinker, beer)
Frequents(drinker, bar)
Serves(bar, beer)

Drinker Category 3

Find drinkers that frequent some bar that serves some beer they like.

$$\{x \mid \exists F \in \text{Frequents} (F.\text{drinker} = x.\text{drinker} \wedge \exists S \in \text{Serves} \exists L \in \text{Likes} (F.\text{drinker} = L.\text{drinker}) \wedge (F.\text{bar} = S.\text{bar}) \wedge (S.\text{beer} = L.\text{beer})))\}$$

Find drinkers that frequent only bars that serve some beer they like.

$$\{x \mid \exists F \in \text{Frequents} (F.\text{drinker} = x.\text{drinker}) \wedge [\forall F1 \in \text{Frequents} (F.\text{drinker} = F1.\text{drinker}) \Rightarrow \exists S \in \text{Serves} \exists L \in \text{Likes} ((F1.\text{bar} = S.\text{bar}) \wedge (F1.\text{drinker} = L.\text{drinker}) \wedge (S.\text{beer} = L.\text{beer}))]\}$$

Find drinkers that frequent some bar that serves only beers they like.

$$\{x \mid \exists F \in \text{Frequents} (F.\text{drinker} = x.\text{drinker}) \wedge [\forall S \in \text{Serves} (F.\text{bar} = S.\text{bar}) \Rightarrow \exists L \in \text{Likes} ((F.\text{drinker} = L.\text{drinker}) \wedge (S.\text{beer} = L.\text{beer}))]\}$$

21

Likes(drinker, beer)
Frequents(drinker, bar)
Serves(bar, beer)

Drinker Category 4

Find drinkers that frequent some bar that serves some beer they like.

$$\{x \mid \exists F \in \text{Frequents} (F.\text{drinker} = x.\text{drinker} \wedge \exists S \in \text{Serves} \exists L \in \text{Likes} (F.\text{drinker} = L.\text{drinker}) \wedge (F.\text{bar} = S.\text{bar}) \wedge (S.\text{beer} = L.\text{beer})))\}$$

Find drinkers that frequent only bars that serve some beer they like.

$$\{x \mid \exists F \in \text{Frequents} (F.\text{drinker} = x.\text{drinker}) \wedge [\forall F1 \in \text{Frequents} (F.\text{drinker} = F1.\text{drinker}) \Rightarrow \exists S \in \text{Serves} \exists L \in \text{Likes} ((F1.\text{bar} = S.\text{bar}) \wedge (F1.\text{drinker} = L.\text{drinker}) \wedge (S.\text{beer} = L.\text{beer}))]\}$$

Find drinkers that frequent some bar that serves only beers they like.

$$\{x \mid \exists F \in \text{Frequents} (F.\text{drinker} = x.\text{drinker}) \wedge [\forall S \in \text{Serves} (F.\text{bar} = S.\text{bar}) \Rightarrow \exists L \in \text{Likes} ((F.\text{drinker} = L.\text{drinker}) \wedge (S.\text{beer} = L.\text{beer}))]\}$$

Find drinkers that frequent only bars that serve only beer they like.

...

22

Likes(drinker, beer)
Frequents(drinker, bar)
Serves(bar, beer)

Drinker Category 4

Find drinkers that frequent some bar that serves some beer they like.

$$\{x \mid \exists F \in \text{Frequents} (F.\text{drinker} = x.\text{drinker} \wedge \exists S \in \text{Serves} \exists L \in \text{Likes} (F.\text{drinker} = L.\text{drinker}) \wedge (F.\text{bar} = S.\text{bar}) \wedge (S.\text{beer} = L.\text{beer})))\}$$

Find drinkers that frequent only bars that serve some beer they like.

$$\{x \mid \exists F \in \text{Frequents} (F.\text{drinker} = x.\text{drinker}) \wedge [\forall F1 \in \text{Frequents} (F.\text{drinker} = F1.\text{drinker}) \Rightarrow \exists S \in \text{Serves} \exists L \in \text{Likes} ((F1.\text{bar} = S.\text{bar}) \wedge (F1.\text{drinker} = L.\text{drinker}) \wedge (S.\text{beer} = L.\text{beer}))]\}$$

Find drinkers that frequent some bar that serves only beers they like.

$$\{x \mid \exists F \in \text{Frequents} (F.\text{drinker} = x.\text{drinker}) \wedge [\forall S \in \text{Serves} (F.\text{bar} = S.\text{bar}) \Rightarrow \exists L \in \text{Likes} ((F.\text{drinker} = L.\text{drinker}) \wedge (S.\text{beer} = L.\text{beer}))]\}$$

Find drinkers that frequent only bars that serve only beer they like.

$$\{x \mid \exists F \in \text{Frequents} (F.\text{drinker} = x.\text{drinker}) \wedge [\forall F1 \in \text{Frequents} (F.\text{drinker} = F1.\text{drinker}) \Rightarrow [\forall S \in \text{Serves} (F1.\text{bar} = S.\text{bar}) \Rightarrow \exists L \in \text{Likes} ((F.\text{drinker} = L.\text{drinker}) \wedge (S.\text{beer} = L.\text{beer}))]\]$$

23

Why should we care about RC

- RC is declarative, like SQL, and unlike RA (which is operational)
- Gives foundation of database queries in first-order logic
 - you cannot express all aggregates in RC, e.g. cardinality of a relation or sum (possible in extended RA and SQL)
 - still can express conditions like “at least two tuples” (or any constant)
- RC expression may be much simpler than SQL queries
 - and easier to check for correctness than SQL
 - power to use \forall and \Rightarrow
 - then you can systematically go to a “correct” SQL or RA query

Duke CS, Fall 2019 CompSci 516: Database Systems 24

Likes(drinker, beer)
Frequents(drinker, bar)
Serves(bar, beer)

Drinker category 5!

From RC to SQL

Query: Find drinkers that like some beer (so much) that they frequent all bars that serve it

$$\{x \mid \exists L \in \text{Likes} (L.\text{drinker} = x.\text{drinker}) \wedge [\forall S \in \text{Serves} (L.\text{beer} = S.\text{beer}) \Rightarrow \exists F \in \text{Frequents} ((F.\text{drinker} = L.\text{drinker}) \wedge (S.\text{beer} = L.\text{beer}))]\}$$

Duke CS, Fall 2019 CompSci 516: Database Systems 25

Likes(drinker, beer)
Frequents(drinker, bar)
Serves(bar, beer)

From RC to SQL (or RA)

Query: Find drinkers that like some beer so much that they frequent all bars that serve it

$$\{x \mid \exists L \in \text{Likes} (L.\text{drinker} = x.\text{drinker}) \wedge [\forall S \in \text{Serves} [(L.\text{beer} = S.\text{beer}) \Rightarrow \exists F \in \text{Frequents} ((F.\text{drinker} = L.\text{drinker}) \wedge (S.\text{beer} = L.\text{beer}))]]\}$$

$$\equiv \{x \mid \exists L \in \text{Likes} (L.\text{drinker} = x.\text{drinker}) \wedge [\forall S \in \text{Serves} [\neg (L.\text{beer} = S.\text{beer}) \vee \exists F \in \text{Frequents} ((F.\text{drinker} = L.\text{drinker}) \wedge (S.\text{beer} = L.\text{beer}))]]\}$$

Step 1: Replace \forall with \exists using de Morgan's Laws

$$Q(x) = \exists y. \text{Likes}(x, y) \wedge [\neg \exists S \in \text{Serves} [(L.\text{beer} = S.\text{beer}) \wedge \neg \exists F \in \text{Frequents} ((F.\text{drinker} = L.\text{drinker}) \wedge (S.\text{beer} = L.\text{beer}))]]$$

$\forall x P(x)$ same as $\neg \exists x \neg P(x)$

SQL or RA does not have \forall !
Now you got all \exists and \neg expressible in RA/SQL

$\neg(\neg \forall x Q)$ same as $P \wedge \neg Q$

Duke CS, Fall 2019 CompSci 516: Database Systems 26

Likes(drinker, beer)
Frequents(drinker, bar)
Serves(bar, beer)

From RC to SQL

Query: Find drinkers that like some beer so much that they frequent all bars that serve it

$$Q(x) = \exists y. \text{Likes}(x, y) \wedge \neg \exists S \in \text{Serves} [(L.\text{beer} = S.\text{beer}) \wedge \neg \exists F \in \text{Frequents} ((F.\text{drinker} = L.\text{drinker}) \wedge (S.\text{beer} = L.\text{beer}))]$$

Step 2: Translate into SQL

```
SELECT DISTINCT L.drinker
FROM Likes L
WHERE not exists
  (SELECT S.bar
   FROM Serves S
   WHERE L.beer=S.beer
   AND not exists (SELECT *
                   FROM Frequents F
                   WHERE F.drinker=L.drinker
                   AND F.bar=S.bar))
```

We will see a "methodical and correct" translation through "safe queries" in Datalog

Duke CS, Fall 2019 CompSci 516: Database Systems 27

Database Normalization

Recap from Lecture-5

ssn (S)	name (N)	lot (L)	rating (R)	hourly-wage (W)	hours-worked (H)
111-11-1111	Attishoo	48	8	10	40
222-22-2222	Smiley	22	8	10	30
333-33-3333	Smethurst	35	5	7	30
444-44-4444	Guldu	35	5	7	32
555-55-5555	Madayan	35	8	10	40

Redundancy is bad! (well...not always?)

- Redundant storage Solution: Decomposition!
- Update anomalies
- Insertion anomalies Be careful about "Lossy decomposition"!
- Deletion anomalies (on blackboard)

Schema is forcing to store (complex) associations among tuples
Nulls may or may not help

Duke CS, Fall 2019 CompSci 516: Database Systems 29

Decompositions should be used judiciously

- Do we need to decompose a relation?
 - Several "normal forms" exist to identify possible redundancy at different granularity
 - If a relation is not in one of them, may need to decompose further
- What are the problems with decomposition?
 - Bad decompositions: e.g., Lossy decompositions
 - Performance issues -- decomposition may both
 - help performance (for updates, some queries accessing part of data), or
 - hurt performance (new joins may be needed for some queries)

Duke CS, Fall 2019 CompSci 516: Database Systems 30

Functional Dependencies (FDs)

- A functional dependency (FD) $X \rightarrow Y$ holds over relation R if, for every allowable instance r of R:
 - i.e., given two tuples in r , if the X values agree, then the Y values must also agree
 - X and Y are sets of attributes
 - $t1 \in r, t2 \in r, \pi_X(t1) = \pi_X(t2)$ implies $\pi_Y(t1) = \pi_Y(t2)$

A	B	C	D
a1	b1	c1	d1
a1	b1	c1	d2
a1	b2	c2	d1
a2	b1	c3	d1

What is a (possible) FD here?

Duke CS, Fall 2019 CompSci 516: Database Systems 31

Functional Dependencies (FDs)

- A functional dependency (FD) $X \rightarrow Y$ holds over relation R if, for every allowable instance r of R:
 - i.e., given two tuples in r , if the X values agree, then the Y values must also agree
 - X and Y are sets of attributes
 - $t1 \in r, t2 \in r, \pi_X(t1) = \pi_X(t2)$ implies $\pi_Y(t1) = \pi_Y(t2)$

A	B	C	D
a1	b1	c1	d1
a1	b1	c1	d2
a1	b2	c2	d1
a2	b1	c3	d1

What is a (possible) FD here?

$AB \rightarrow C$

Note that, AB is not a key

Duke CS, Fall 2019 CompSci 516: Database Systems 32

Can we detect FDs from an instance?

- An FD is a statement about **all** allowable relation instances
 - Must be identified based on semantics of application
 - Given some allowable instance $r1$ of R, we can check if it **violates** some FD f , but we **cannot tell if f holds over R**
- K is a candidate key for R means that $K \rightarrow R$
 - denoting R = all attributes of R too
 - However, $S \rightarrow R$ does not require S to be minimal
 - e.g. S can be a superkey

Duke CS, Fall 2019 CompSci 516: Database Systems 33

FD from a key

- Consider a relation R(A,B, C, D) where AB is a key
- Which FD must hold on R?

- $AB \rightarrow ABCD$

- However, $S \rightarrow ABCD$ does not mean S is a key. Why?
 - S can be a superkey!
 - E.g., $ABC \rightarrow ABCD$ in R, but ABC is not a key

Duke CS, Fall 2019 CompSci 516: Database Systems 34

Armstrong's Axioms

- X, Y, Z are sets of attributes

A	B	C	D
a1	b1	c1	d1
a1	b1	c1	d2
a1	b2	c2	d1
a2	b1	c3	d1

Apply these rules on $AB \rightarrow C$ and check

- Reflexivity:** If $X \supseteq Y$, then $X \rightarrow Y$, e.g., $ABC \rightarrow AB$
- Augmentation:** If $X \rightarrow Y$, then $XZ \rightarrow YZ$ for any Z,
 - e.g., $AB \rightarrow C \Rightarrow ABDE \rightarrow CDE$
- Transitivity:** If $X \rightarrow Y$ and $Y \rightarrow Z$, then $X \rightarrow Z$
 - e.g., $AB \rightarrow C$ and $C \rightarrow D \Rightarrow AB \rightarrow D$

- Additional rules that follow from Armstrong's Axioms
 - $A \rightarrow B$ and $A \rightarrow C \Rightarrow A \rightarrow BC$
 - $A \rightarrow BC \Rightarrow A \rightarrow B, A \rightarrow C$

- Union:** If $X \rightarrow Y$ and $X \rightarrow Z$, then $X \rightarrow YZ$
- Decomposition:** If $X \rightarrow YZ$, then $X \rightarrow Y$ and $X \rightarrow Z$

Duke CS, Fall 2019 CompSci 516: Database Systems 35

Closure of a set of FDs

- Given some FDs, we can usually infer additional FDs:
 - $SSN \rightarrow DEPT$, and $DEPT \rightarrow LOT$ implies $SSN \rightarrow LOT$
- An FD f is **implied** by a set of FDs F if f holds whenever all FDs in F hold.
- F^* = closure of FDs F is the set of all FDs that are implied by F
- S^* = closure of attributes S is the set of all attributes that are implied by S according to F^*

Armstrong's Axioms are **sound** and **complete** inference rules for FDs

- sound: they only generate FDs in closure F^* for F
- complete: by repeated application of these rules, all FDs in F^* will be generated

Duke CS, Fall 2019 CompSci 516: Database Systems 36

Computing Attribute Closure

Algorithm:

Let's do the example first, Then look at the algo yourself

- closure = X
- Repeat until no change
 - if there is an FD $U \rightarrow V$ in F such that $U \subseteq \text{closure}$, then $\text{closure} = \text{closure} \cup V$

Does $F = \{A \rightarrow B, B \rightarrow C, CD \rightarrow E\}$ imply

1. $A \rightarrow E$? (i.e. is $A \rightarrow E$ in the closure F^+ , or E in A^+ ?)
2. $AD \rightarrow E$?

On blackboard

Duke CS, Fall 2019
CompSci 516: Database Systems
37

Normal Forms

- Question: given a schema, how to decide whether any schema refinement is needed at all?
- If a relation is in a certain normal forms, it is known that certain kinds of problems are avoided/minimized
- Helps us decide whether decomposing the relation is something we want to do

Duke CS, Fall 2019
CompSci 516: Database Systems
38

FDs play a role in detecting redundancy

Example

- Consider a relation R with 3 attributes, ABC
 - No FDs hold: There is no redundancy here – no decomposition needed
 - Given $A \rightarrow B$: Several tuples could have the same A value, and if so, they'll all have the same B value \Rightarrow redundancy \Rightarrow decomposition may be needed if A is not a key
- Intuitive idea:
 - if there is any non-key dependency, e.g. $A \rightarrow B$, decompose!

Duke CS, Fall 2019
CompSci 516: Database Systems
39

Normal Forms

- R is in 4NF
- \Rightarrow R is in BCNF
- \Rightarrow R is in 3NF
- \Rightarrow R is in 2NF (a historical one)
- \Rightarrow R is in 1NF (every field has atomic values)

Only BCNF and 4NF are covered in the class

Duke CS, Fall 2019
CompSci 516: Database Systems
40

Boyce-Codd Normal Form (BCNF)

- Relation R with FDs F is in BCNF if, for all $X \rightarrow A$ in F
 - $A \in X$ (called a trivial FD), or
 - X contains a key for R
 - i.e. X is a superkey

Duke CS, Fall 2019
CompSci 516: Database Systems
41

BCNF decomposition algorithm

- Find a BCNF violation
 - That is, a non-trivial FD $X \rightarrow Y$ in R where X is not a super key of R
- Decompose R into R_1 and R_2 , where
 - R_1 has attributes $X \cup Y$
 - R_2 has attributes $X \cup Z$, where Z contains all attributes of R that are in neither X nor Y
- Repeat until all relations are in BCNF
- Also gives a lossless decomposition!
 - Check yourself

Duke CS, Fall 2019
CompSci 516: Database Systems
42

BCNF decomposition example - 1

On blackboard

- $CSJDPQV$, key C , $F = \{JP \rightarrow C, SD \rightarrow P, J \rightarrow S\}$
 - To deal with $SD \rightarrow P$, decompose into SDP , $CSJDQV$.
 - To deal with $J \rightarrow S$, decompose $CSJDQV$ into JS and $CJDQV$
- Is $JP \rightarrow C$ a violation of BCNF?
- Note:
 - several dependencies may cause violation of BCNF
 - The order in which we pick them may lead to very different sets of relations
 - there may be multiple correct decompositions (can pick $J \rightarrow S$ first)

Duke CS, Fall 2019 CompSci 516: Database Systems 43

BCNF decomposition example - 2

$uid \rightarrow uname, twitterid$
 $twitterid \rightarrow uid$
 $uid, gid \rightarrow fromDate$

$UserJoinsGroup (uid, uname, twitterid, gid, fromDate)$
BCNF violation: $uid \rightarrow uname, twitterid$

$User (uid, uname, twitterid)$
 $uid \rightarrow uname, twitterid$
 $twitterid \rightarrow uid$
BCNF

$Member (uid, gid, fromDate)$
 $uid, gid \rightarrow fromDate$
BCNF

Duke CS, Fall 2019 CompSci 516: Database Systems 44

BCNF decomposition example - 3

It is not enough to only look at given FDs! You need to consider the closure!

$uid \rightarrow uname, twitterid$
 $twitterid \rightarrow uid$
 $uid, gid \rightarrow fromDate$

$UserJoinsGroup (uid, uname, twitterid, gid, fromDate)$
BCNF violation: $twitterid \rightarrow uid$

$UserId (twitterid, uid)$
BCNF

$UserJoinsGroup' (twitterid, uname, gid, fromDate)$
 $twitterid \rightarrow uname$
 $twitterid, gid \rightarrow fromDate$
BCNF violation: $twitterid \rightarrow uname$

$UserName (twitterid, uname)$
BCNF

$Member (twitterid, gid, fromDate)$
BCNF

apply Armstrong's axioms and rules!

Duke CS, Fall 2019 CompSci 516: Database Systems 45

Recap

- Functional dependencies: a generalization of the key concept
- Non-key functional dependencies: a source of redundancy
- BCNF decomposition: a method for removing redundancies
 - BCNF decomposition is a lossless join decomposition
- BCNF: schema in this normal form has no redundancy due to FD's

Duke CS, Fall 2019 CompSci 516: Database Systems 46

BCNF = no redundancy?

- $User (uid, gid, place)$
 - A user can belong to multiple groups
 - A user can register places she's visited
 - Groups and places have nothing to do with other
 - FD's?
 - None
 - BCNF?
 - Yes
 - Redundancies?
 - Tons!

uid	gid	place
142	dps	Springfield
142	dps	Australia
456	abc	Springfield
456	abc	Morocco
456	gov	Springfield
456	gov	Morocco
...

Duke CS, Fall 2019 CompSci 516: Database Systems 47

Multivalued dependencies

- A multivalued dependency (MVD) has the form $X \twoheadrightarrow Y$, where X and Y are sets of attributes in a relation R
- $X \twoheadrightarrow Y$ means that whenever two rows in R agree on all the attributes of X , then we can swap their Y components and get two rows that are also in R

X	Y	Z
a	b ₁	c ₁
a	b ₂	c ₂
a	b ₂	c ₁
a	b ₁	c ₂
...

Duke CS, Fall 2019 CompSci 516: Database Systems 48

MVD examples

User (*uid, gid, place*)

- $uid \twoheadrightarrow gid$
- $uid \twoheadrightarrow place$
 - Intuition: given *uid*, attributes *gid* and *place* are "independent"
- $uid, gid \twoheadrightarrow place$
 - Trivial: $LHS \cup RHS = \text{all attributes of } R$
- $uid, gid \twoheadrightarrow uid$
 - Trivial: $LHS \supseteq RHS$

Duke CS, Fall 2019 CompSci 516: Database Systems 49

Read this slide after looking at the examples

An elegant solution: "chase"

- Given a set of FD's and MVD's \mathcal{D} , does another dependency d (FD or MVD) follow from \mathcal{D} ?
- Procedure
 - Start with the premise of d , and treat them as "seed" tuples in a relation
 - Apply the given dependencies in \mathcal{D} repeatedly
 - If we apply an FD, we infer equality of two symbols
 - If we apply an MVD, we infer more tuples
 - If we infer the conclusion of d , we have a **proof**
 - Otherwise, if nothing more can be inferred, we have a **counterexample**

Duke CS, Fall 2019 CompSci 516: Database Systems 50

Proof by chase

- In $R(A, B, C, D)$, does $A \twoheadrightarrow B$ and $B \twoheadrightarrow C$ imply that $A \twoheadrightarrow C$?

Have:

A	B	C	D
a	b ₁	c ₁	d ₁
a	b ₂	c ₂	d ₂

 Need:

A	B	C	D
a	b ₁	c ₂	d ₁
a	b ₂	c ₁	d ₂

$A \twoheadrightarrow B$

A	B	C	D
a	b ₂	c ₁	d ₁
a	b ₁	c ₂	d ₂

$B \twoheadrightarrow C$

A	B	C	D
a	b ₂	c ₂	d ₁
a	b ₁	c ₁	d ₂

$B \twoheadrightarrow C$

A	B	C	D
a	b ₁	c ₂	d ₁
a	b ₂	c ₁	d ₂

Duke CS, Fall 2019 CompSci 516: Database Systems 51

Another proof by chase

- In $R(A, B, C, D)$, does $A \twoheadrightarrow B$ and $B \twoheadrightarrow C$ imply that $A \twoheadrightarrow C$?

Have:

A	B	C	D
a	b ₁	c ₁	d ₁
a	b ₂	c ₂	d ₂

 Need: $c_1 = c_2$

$A \twoheadrightarrow B$ $b_1 = b_2$

$B \twoheadrightarrow C$ $c_1 = c_2$

In general, with both MVD's and FD's, chase can generate both new tuples and new equalities

Duke CS, Fall 2019 CompSci 516: Database Systems 52

Counterexample by chase

- In $R(A, B, C, D)$, does $A \twoheadrightarrow BC$ and $CD \twoheadrightarrow B$ imply that $A \twoheadrightarrow B$?

Have:

A	B	C	D
a	b ₁	c ₁	d ₁
a	b ₂	c ₂	d ₂

 Need: $b_1 = b_2$

$A \twoheadrightarrow BC$

A	B	C	D
a	b ₂	c ₂	d ₁
a	b ₁	c ₁	d ₂

Counterexample!

Duke CS, Fall 2019 CompSci 516: Database Systems 53

4NF

- A relation R is in **Fourth Normal Form (4NF)** if
 - For every non-trivial MVD $X \twoheadrightarrow Y$ in R , X is a **superkey**
 - That is, all FD's and MVD's follow from "key \rightarrow other attributes" (i.e., no MVD's and no FD's besides key functional dependencies)
- 4NF is stronger than BCNF
 - Because every FD is also a MVD

Duke CS, Fall 2019 CompSci 516: Database Systems 54

4NF decomposition algorithm

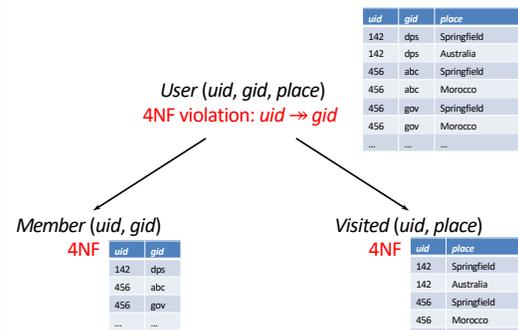
- Find a 4NF violation
 - A non-trivial MVD $X \twoheadrightarrow Y$ in R where X is not a superkey
- Decompose R into R_1 and R_2 , where
 - R_1 has attributes $X \cup Y$
 - R_2 has attributes $X \cup Z$ (where Z contains R attributes not in X or Y)
- Repeat until all relations are in 4NF
- Almost identical to BCNF decomposition algorithm
- Any decomposition on a 4NF violation is lossless

Duke CS, Fall 2019

CompSci 516: Database Systems

55

4NF decomposition example



Duke CS, Fall 2019

CompSci 516: Database Systems

56

Other kinds of dependencies and normal forms

- Dependency preserving decompositions
- Join dependencies
- Inclusion dependencies
- 5NF, 3NF, 2NF
- See book if interested (not covered in class)

Duke CS, Fall 2019

CompSci 516: Database Systems

57

Summary

- Philosophy behind BCNF, 4NF:
 - Data should depend on the key, the whole key, and nothing but the key!
 - You could have multiple keys though
- Redundancy is not desired typically
 - not always, mainly due to performance reasons
- Functional/multivalued dependencies – capture redundancy
- Decompositions – eliminate dependencies (should not be lossy!)
- Normal forms
 - Guarantees certain non-redundancy
 - BCNF, and 4NF
- How to decompose into BCNF, 4NF
- Chase



Duke CS, Fall 2019

CompSci 516: Database Systems

58