

## G53RDB: Theory of Relational Databases Lecture 13

Natasha Alechina  
School of Computer Science & IT  
nza@cs.nott.ac.uk

### Plan of the lecture

- Relational calculus and relational algebra
- Deductive databases
- Datalog:
  - syntax
  - non-recursive queries
  - meaning of non-recursive queries

Lecture 16

2

### Formulas and relations

- A formula with n free variables defines an n-ary relation: the set of n-tuples which satisfy this formula.
- For example,  $\exists x \text{ ChildOf}(x,y,z)$  defines a binary relation  $\{ \langle \mathbf{a}, \mathbf{b} \rangle : \langle y:\mathbf{a},z:\mathbf{b} \rangle \text{ satisfy } \exists z \text{ ChildOf}(x,y,z) \}$

ChildOf(x,y,z)

x	y	z
Mary	Tom	Jane
Bill	Tom	Jane

$\exists x \text{ ChildOf}(x,y,z)$

y	z
Tom	Jane

Lecture 16

3

### Relational calculus and relational algebra

The relationship is really quite simple (if we abstract from many technical details like variable names vs attributes):

- Starting from the same set of relations, all relations we can build using relational algebra operators, we can define by a formula of relational calculus, and vice versa.
- The two formalisms are equally expressive: every query we can define in one, we can define in another.
- Certain things are easier with relational calculus: it is more flexible, closer to natural language. On the other hand, it is much easier to design algorithms for evaluating algebraic expressions.

Lecture 16

4

### Relations and formulas

- Suppose we have a set of relations (database tables). We can define new relations using relational algebra operators  $\pi, \sigma, \times, \cup, -$ .
- I am going to show that for each such relation, we can write a first order formula which is true exactly for the set of tuples in the relation.

Lecture 16

5

### Basic relations and formulas

- For each basic n-ary relation  $\mathbf{R}$  (a database table), we introduce a relation name  $R$
- The formula  $R(x_1, \dots, x_n)$  is then true for exactly tuples in  $\mathbf{R}$ .
- Note that instead of  $x_1, \dots, x_n$ , we could have used  $y_1, \dots, y_n$ .

Relation  $\mathbf{R}$

$\langle \mathbf{a}_1, \dots, \mathbf{a}_n \rangle$
$\langle \mathbf{b}_1, \dots, \mathbf{b}_n \rangle$
$\langle \mathbf{c}_1, \dots, \mathbf{c}_n \rangle$

Relation name  $R$

$R(x_1, \dots, x_n)$

Lecture 16

6

## Inductive step

- The rest we show by induction.
- We can define basic relations by atomic formulas, so now need to show how we can define union, difference, product, projection and selection of two definable relations.

Lecture 16

7

## Union of relations

- Suppose that we have two union compatible relations  $\mathbf{R}$  and  $\mathbf{P}$ .
- Inductive hypothesis: there exists a first order formula  $\phi(x_1, \dots, x_n)$  which defines relation  $\mathbf{R}$  and a first order formula  $\psi(x_1, \dots, x_n)$  which defines relation  $\mathbf{P}$  (a tuple satisfies the formula if, and only if, it belongs to the relation).
- Then  $\mathbf{R} \cup \mathbf{P}$  is defined by  $\phi(x_1, \dots, x_n) \vee \psi(x_1, \dots, x_n)$ .
- Indeed, a tuple is in  $\mathbf{R} \cup \mathbf{P}$  if, and only if, it belongs to  $\mathbf{R}$  or  $\mathbf{P}$ ; by the inductive hypothesis, this holds if, and only if, it satisfies  $\phi(x_1, \dots, x_n)$  or  $\psi(x_1, \dots, x_n)$ ; which means, if, and only if, it satisfies  $\phi(x_1, \dots, x_n) \vee \psi(x_1, \dots, x_n)$ .

Lecture 16

8

## Difference of relations

- Suppose that we have two union compatible relations  $\mathbf{R}$  and  $\mathbf{P}$ .
- Inductive hypothesis: there exists a first order formula  $\phi(x_1, \dots, x_n)$  which defines relation  $\mathbf{R}$  and a first order formula  $\psi(x_1, \dots, x_n)$  which defines relation  $\mathbf{P}$ .
- Then  $\mathbf{R} - \mathbf{P}$  is defined by  $\phi(x_1, \dots, x_n) \& \neg\psi(x_1, \dots, x_n)$ .
- Indeed, a tuple is in  $\mathbf{R} - \mathbf{P}$  if, and only if, it belongs to  $\mathbf{R}$  and does not belong to  $\mathbf{P}$ ; by the inductive hypothesis, this holds if, and only if, it satisfies  $\phi(x_1, \dots, x_n)$  and does not satisfy  $\psi(x_1, \dots, x_n)$ ; which means, if, and only if, it satisfies  $\phi(x_1, \dots, x_n) \& \neg\psi(x_1, \dots, x_n)$ .

Lecture 16

9

## Product of relations

- Suppose that we have two relations  $\mathbf{R}$  and  $\mathbf{P}$ .
- Inductive hypothesis: there exists a first order formula  $\phi(x_1, \dots, x_n)$  which defines relation  $\mathbf{R}$  and a first order formula  $\psi(y_1, \dots, y_k)$  which defines relation  $\mathbf{P}$ .
- Then  $\mathbf{R} \times \mathbf{P}$  is defined by  $\phi(x_1, \dots, x_n) \& \psi(y_1, \dots, y_k)$ .
- Indeed, a tuple  $\langle \mathbf{a}_1, \dots, \mathbf{a}_n, \mathbf{b}_1, \dots, \mathbf{b}_k \rangle$  is in  $\mathbf{R} \times \mathbf{P}$  if, and only if, first  $n$  values belong to  $\mathbf{R}$  and next  $k$  values belong to  $\mathbf{P}$ ; by the inductive hypothesis, this holds if, and only if,  $\langle \mathbf{a}_1, \dots, \mathbf{a}_n \rangle$  satisfies  $\phi(x_1, \dots, x_n)$  and  $\langle \mathbf{b}_1, \dots, \mathbf{b}_k \rangle$  satisfies  $\psi(y_1, \dots, y_k)$ ; which means, if, and only if, it satisfies  $\phi(x_1, \dots, x_n) \& \psi(y_1, \dots, y_k)$ .

Lecture 16

10

## Projection

- Suppose we have a relation  $\mathbf{R}$ , and by the inductive hypothesis, there exists a first order formula  $\phi(x_1, \dots, x_n, y_1, \dots, y_k)$  which defines  $\mathbf{R}$ . Suppose we want to project just on the attribute positions corresponding to  $x$ s.
- Then  $\pi_{x_1, \dots, x_n} \mathbf{R}$  is defined by  $\exists y_1 \dots \exists y_n \phi$ :

Relation  $\mathbf{R}$

$x_1$	$\dots$	$x_n$	$y_1$	$\dots$	$y_k$
$\mathbf{a}_1$	$\dots$	$\mathbf{a}_n$	$\mathbf{b}_1$	$\dots$	$\mathbf{b}_k$
$\mathbf{c}_1$	$\dots$	$\mathbf{c}_n$	$\mathbf{d}_1$	$\dots$	$\mathbf{d}_k$
$\mathbf{e}_1$	$\dots$	$\mathbf{e}_n$	$\mathbf{f}_1$	$\dots$	$\mathbf{f}_k$

Lecture 16

11

## Selection

- Suppose we have a relation  $\mathbf{R}$ , and by the inductive hypothesis, there exists a first order formula  $\phi(x_1, \dots, x_n)$  which defines  $\mathbf{R}$ . Suppose we want to select tuples which satisfy property  $\alpha$ .
- If  $\alpha$  can be expressed in first order logic, as some condition  $\alpha(x_1, \dots, x_n)$ , then
- $\sigma_{\alpha} \mathbf{R}$  is defined by  $\alpha(x_1, \dots, x_n) \& \phi(x_1, \dots, x_n)$ .
- All we need to express  $\alpha$  is to add to our first order language comparison operators and constants for database values. We already have  $\&$ ,  $\neg$  and  $\vee$ .

Lecture 16

12

## So...

- Every relation definable in relational algebra, is also definable in first order logic.
- The opposite direction (everything definable in first order logic, is definable in relational algebra) also holds, but is a lot harder to prove.
- The two formalisms (relational algebra and first order logic) have the same expressive power.

Lecture 16

13

## Deductive databases

- Deductive databases are very similar to relational databases, but their query languages have an additional capability: they can express recursive queries.
- The syntax is closely based on predicate logic

Lecture 16

14

## Datalog

- One of the query languages used in deductive databases is called Datalog.
- Its syntax is very similar to predicate logic/relational calculus
- It is also very similar to logical programming (Prolog), but there are some important differences.

Lecture 16

15

## Syntax of Datalog

- Similarly to relational calculus, Datalog statements are built from atoms.
- **Relational atom:** if R is a relation name of arity n and  $t_1, \dots, t_n$  are variables or constants, then  $R(t_1, \dots, t_n)$  is a relational atom
- **Arithmetic atom:** comparison between two arithmetic expressions, for example  $x > 5$ ,  $x \geq y$ , etc.

Lecture 16

16

## Datalog syntax: rules

- A Datalog **rule** is an expression of the form
$$R_1 \leftarrow R_2 \text{ AND } \dots \text{ AND } R_n$$
where  $n \geq 1$ ,  $R_1$  is a relational atom, and  $R_2, \dots, R_n$  are relational or arithmetic atoms, possibly preceded by NOT.
- $R_1$  is called the **head** of the rule and  $R_2, \dots, R_n$  the **body** of the rule.
- $R_2, \dots, R_n$  are called **subgoals**.

Lecture 16

17

## Example

- Suppose we have a relation Person over schema (Name, Age, Address, Telephone). Then the following Datalog rule will define a relation which contains names of people aged over 18:

$$\text{Adult}(x) \leftarrow \text{Person}(x,y,z,u) \text{ AND } y \geq 18$$

Lecture 16

18

## Anonymous variables

$Adult(x) \leftarrow Person(x,y,z,u) \text{ AND } y \geq 18$

- Note that z and u appear only once in the rule, and their names do not matter (they are not compared to any other variable and do not appear in the definition of Adult).
- In such cases (so that a programmer does not have to invent new names for dummy variables) a variable may be replaced by an underscore:

$Adult(x) \leftarrow Person(x,y,\_,\_) \text{ AND } y \geq 18$

Lecture 16

19

## Datalog query

- A **Datalog query** is a finite set of Datalog rules
- If there is only one relation which appears as a head of a rule in the query, the tuples in that relation are taken as the answer to the query.
- For example,

$Parent(x,y) \leftarrow Mother(x,y)$

$Parent(x,y) \leftarrow Father(x,y)$

defines Parent relation (using relations Father and Mother)

- If there is more than one relation appearing as a head, one of them is the main predicate to be defined and others are auxiliary.

Lecture 16

20

## Datalog and logic programming (Prolog)

- Datalog programs look very similar to logic programs.
- Main differences are:
  - there are no functional symbols in Datalog
  - the use of programs is different:
    - Datalog program is a mapping from the database to the set of new relations; it is assumed that the program is small and the database is large
    - Logic program contains (a small amount of) data as part of the program; it would just list “facts” such as  $Mother(Jane, Mary) \leftarrow$

Lecture 16

21

## Meaning of Datalog rules

- First approximation (non-recursive queries):
  - take the values of variables which make the body of the rule true (make each subgoal true; NOT R is true if R is false)
  - see what values the variables of the head take;
  - add the resulting tuple to the predicate in the head of the rule.

Lecture 16

22

## Example

- Suppose we have a relation Person over schema (Name, Age, Address, Telephone). Then the following Datalog rule will define a relation which contains names of people aged over 18:

$Adult(x) \leftarrow Person(x,y,z,u) \text{ AND } y \geq 18$

- We take all  $\langle \text{name, age, addr, tel} \rangle$  in Person for which it is also true that  $\text{age} \geq 18$ , and add  $\langle \text{name} \rangle$  to Adult.
- Adult relation is the same as  $\pi_{\text{Name}} \sigma_{\text{Age} \geq 18}(\text{Person})$ .

Lecture 16

23

## Another example

$Parent(x,y) \leftarrow Mother(x,y)$

$Parent(x,y) \leftarrow Father(x,y)$

- suppose Mother contains tuples  $\langle \text{Jane, Mary} \rangle$  and  $\langle \text{Jane, Bill} \rangle$  and Father contains tuples  $\langle \text{Tom, Mary} \rangle$ ,  $\langle \text{Tom, Bill} \rangle$ .
- All four tuples are added to Parent.

Lecture 16

24

## Example with negation

- Suppose we have a relation Person over schema (Name, Age, Address, Telephone).  
$$\text{Child}(x) \leftarrow \text{Person}(x,y,z,u) \text{ AND NOT}(y \geq 18)$$
- We take all  $\langle \text{name, age, addr, tel} \rangle$  in Person for which it is also true that  $\text{NOT}(\text{age} \geq 18)$ , and add  $\langle \text{name} \rangle$  to Adult.
- $\text{NOT}(\text{age} \geq 18)$  is true if  $\text{age} \geq 18$  is false, so we add all tuples where  $\text{age} < 18$ .

Lecture 16

25

## Safe queries

- We want the result of a query to be a finite relation.
- To ensure this, the following *safety condition* is required:  
*every variable that appears anywhere in the rule must appear in some non-negated relational subgoal.*
- The reason for this is that infinitely many values may satisfy an arithmetical subgoal (e.g.  $x > 0$ ) and infinitely many values are NOT in some finite table of a relation R.

Lecture 16

26

## Datalog and relational algebra

- Every relation definable in relational algebra is definable in Datalog.
- Again we assume that we have a relational name (predicate symbol) R for every basic relation R.
- Then for every operation of relational algebra, we show how to write a corresponding Datalog query.

Lecture 16

27

## Union

- Union of R and S:

$$U(x_1, \dots, x_n) \leftarrow R(x_1, \dots, x_n)$$

$$U(x_1, \dots, x_n) \leftarrow S(x_1, \dots, x_n)$$

Lecture 16

28

## Difference

- Difference of R and S:

$$D(x_1, \dots, x_n) \leftarrow R(x_1, \dots, x_n) \text{ AND NOT } S(x_1, \dots, x_n)$$

Lecture 16

29

## Product

- Product of R and S:

$$P(x_1, \dots, x_n, y_1, \dots, y_k) \leftarrow R(x_1, \dots, x_n) \text{ AND } S(y_1, \dots, y_k)$$

Lecture 16

30

## Projection

- Suppose we want to project  $R$  on attributes  $x_1, \dots, x_n$ .

$$P(x_1, \dots, x_n) \leftarrow R(x_1, \dots, x_n, y_1, \dots, y_k)$$

or

$$P(x_1, \dots, x_n) \leftarrow R(x_1, \dots, x_n, \_ , \dots, \_)$$

Lecture 16

31

## Selection

- Simple case: all conditions in the selection are connected by AND, for example  $\sigma_{\text{Age} > 18 \text{ AND Address} = \text{"London"}}(\text{Person})$

$$\text{Answer}(x, y, z, u) \leftarrow \text{Person}(x, y, z, u) \text{ AND } y > 18 \text{ AND } z = \text{"London"}$$

- If conditions are connected with OR, need more than one rule. For example,  $\sigma_{\text{Age} > 18 \text{ OR Address} = \text{"London"}}(\text{Person})$

$$\text{Answer}(x, y, z, u) \leftarrow \text{Person}(x, y, z, u) \text{ AND } y > 18$$

$$\text{Answer}(x, y, z, u) \leftarrow \text{Person}(x, y, z, u) \text{ AND } z = \text{"London"}$$

Lecture 16

32

## Compound queries

- To translate an arbitrary algebraic expression, create a new predicate for every node in the query tree.
- For example, to do  $\sigma_{\text{Name1} = \text{Name2}}(R \times P)$ :
  - Define predicate  $S = R \times P$
  - Define  $\sigma_{\text{Name1} = \text{Name2}}(S)$

Lecture 16

33

## Informal coursework

- A database of fictitious company contains three relations:
  - GOODS over schema {Producer, ProductCode, Description}
  - DELIVERY over schema {Producer, ProductCode, Branch#, Stock#}
  - STOCK over schema {Branch#, Stock#, Size, Colour, SellPrice, CostPrice, DateIn, DateOut}.

Lecture 16

34

## Define in Datalog

- Query 1: find all producers who supply goods.
- Query 2: find all producers who have delivered goods to any branch of the company.
- Query 3: find SellPrice and CostPrice of all goods delivered to branch L1 still in stock (here, L1 is a value in the attribute domain of Branch#, and products in stock have value InStock for the DateOut attribute).
- Query 4: find Producer, ProductCode, Description for all goods sold at the same day they arrived at any branch.
- Query 5: find Branch#, Size, Colour, SellPrice for all dresses which have not yet been sold (dress is a value in the attribute domain of Description).

Lecture 16

35

## Reading

- Ullman, Widom, chapter 10, or
- Abiteboul, Hull, Vianu chapter 12.

Lecture 16

36