# **Introduction to Databases**

(INFR10080)

(Deductive Databases)

Instructor: Yang Cao

(Fall 2025)



Changelog

v25.0 Initial version

# Datalog rules (1)

$$\underbrace{H(\bar{x})}_{\text{head}} :- \underbrace{S_1(\bar{x}_1), \dots, S_n(\bar{x}_n)}_{\text{body}}$$

Head predicate over variables and constants

Body conjunction of possibly negated atoms (subgoals)

- relational atoms (predicates)
- comparisons between variables/constants
- Head variables are implicitly universally quantified
- Body variables not in head are existentially quantified
- Rule: Body → Head (if the body is true, the head is true)

2/23

## Datalog rules (2)

#### Example

$$Lucky(x) := Customer(x, y, z), Account(u, z, x, w), w > 10000$$

In relational calculus we could write:

Lucky = 
$$\{ x \mid \exists y, z, u, w \text{ Customer}(x, y, z) \land \text{Account}(u, z, x, w) \land w > 10000 \}$$

#### Safety

Every variable (in head or body) appears in at least one non-negated relational atom

**Not safe:** BigNumber(x) :- x > 1000000000

## Datalog programs

#### Program = set of Datalog rules

#### Example

$$Parent(x, y) := Mother(x, y)$$

$$Parent(x, y) := Father(x, y)$$

edb (extensional database) relations stored in the database

can appear only in the body of rules

idb (intensional database) derived relations

can appear both in the head or body of rules

4/23

# From relational algebra to Datalog

Every relational algebra expression can be translated into Datalog

Projection 
$$\pi_{#2}(R)$$

$$E(x) := R(y, x)$$

Selection 
$$\sigma_{\#1 \text{ op } c}(R)$$

$$E(x, y) := R(x, y), x \operatorname{op} c$$

Product 
$$R \times S$$

$$E(x, y, w, z) := R(x, y), S(w, z)$$

Difference 
$$R - S$$

$$E(x, y) := R(x, y), \neg S(x, y)$$

Union 
$$R \cup S$$

$$E(x, y) := R(x, y)$$

$$E(x, y) := S(x, y)$$

## From relational algebra to Datalog

Let *R* and *S* be relations over *A*, *B* 

$$E = \pi_{A,D} \left( \underbrace{\sigma_{B=C} \left( \underbrace{R \times \rho_{A \to C, B \to D}(S)}_{E_1} \right)}_{E_2} \cup \underbrace{\rho_{B \to D} (R - S)}_{E_4} \right)$$

$$E_{1}(x, u, w, y) := R(x, u), S(w, y)$$

$$E_{2}(x, u, w, y) := E_{1}(x, u, w, y), u = w$$

$$E_{3}(x, y) := E_{2}(x, u, w, y)$$

$$E_{4}(x, y) := R(x, y), \neg S(x, y)$$

$$E(x, y) := E_{3}(x, y)$$

$$E(x, y) := E_{4}(x, y)$$

6/23

# Limitations of relational algebra/calculus

Parent = table of pairs x, y where x is the parent of y

$$Parent = \{x, y \mid Parent(x, y)\}$$

$$Grandparent = \{x, y \mid \exists z \ Parent(x, z) \land Parent(z, y)\}$$

$$Great-grandparent = \{x, y \mid \exists z \ Grandparent(x, z) \land Parent(z, y)\}$$

For a given k, we can express the query Ancestor<sup>k</sup>

#### But

We **cannot express** the Ancestor relation itself that is: an Ancestor $^k$  query that works for **every** k

## Recursion in Datalog

The head relation of a rule can appear in its body

```
Ancestor(x, y) := Parent(x, y)
Ancestor(x, y) := Ancestor(x, z), Parent(z, y)
```

#### Intuition

```
x is an ancestor of y if
x is a parent of y
or
x is an ancestor of a parent of y
```

8/23

# Dependency graph

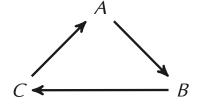
IDB predicate *P* depends on (IDB) predicate *Q* if there is a rule with *P* in the head and *Q* in a subgoal

#### Dependency graph

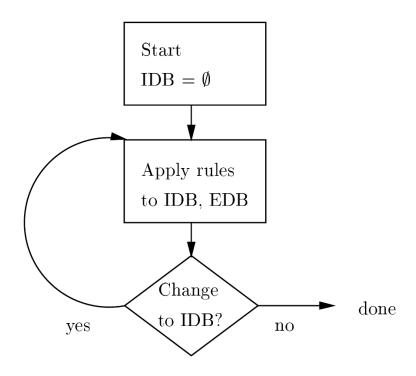
```
nodes IDB predicates edges P \rightarrow Q if P depends on Q
```

A cycle in the dependency graph means the program is recursive

$$C(x) := A(y, x)$$
  
 $C(x) := S(x, y), y > 1$   
 $B(x, y) := C(x), P(x, y)$   
 $A(x, y) := B(y, x)$ 



## **Iterative Fixpoint Evaluation**



10/23

# Evaluation of recursive programs

- The Parent relation (EDB) never changes
- The **Ancestor** relation (IDB) is initially empty

$$\mathsf{Ancestor}_0 = \varnothing$$

• At step i + 1 compute:

Ancestor<sub>$$i+1$$</sub> $(x, y) := Parent(x, y)$   
Ancestor <sub>$i+1$</sub>  $(x, y) := Ancestor $i$  $(x, z)$ , Parent $(z, y)$$ 

Stop when a fixpoint is reached

$$Ancestor_{i+1} = Ancestor_i$$

## Evaluation of recursive programs

	Ancestor	
	John	Mary
IDB:	John	Jane
	Jane	Louis
	Mary	Linda
	Louis	Mark
	John	Linda
	John	Louis
	Jane	Mark
	John	Mark

Par	Parent		
John	Mary		
John	Jane		
Jane	Louis		
Mary	Linda		
Louis	Mark		

EDB:

```
Ancestor(x, y) := Parent(x, y)
Ancestor(x, y) := Ancestor(x, z), Parent(z, y)
```

12/23

## Recursion in SQL

Suppose we have a table Parent with attributes name, child

The definition mimics the structure of the Datalog program

#### Nonlinear recursion

The head relation can appear more than once in its body

```
Ancestor(x, y) :— Parent(x, y)
Ancestor(x, y) :— Ancestor(x, z), Ancestor(x, y)
```

#### Intuition

```
x is an ancestor of y if
x is a parent of y
or
x is an ancestor of an ancestor of y
```

14/23

# Nonlinear recursion in SQL

Not supported

ERROR: recursive reference to query "ancestor" must not appear more than once

## Recursive programs with negation

Consider the program  $P = \{R(x) : -S(x), \neg R(x)\}$ 

Step 0 IDB relation 
$$R$$
 is empty  
Step 1  $R = \{1, 2\}$   
Step 2  $R = \emptyset$   
Step 3  $R = \{1, 2\}$   
Step 4  $R = \emptyset$ 

... No fixpoint!

Iteration never ends

16/23

#### Stratification

Partition a program P into a sequence of subprograms  $P_1, \ldots, P_n$ 

- Each subprogram defines one or more IDB relations
- If a relation *S* is used positively in the definition of *R* then *S* must be defined **earlier or simultaneously** with *R*
- If a relation *S* is used negatively in the definition of *R* then *S* must be defined **strictly before** *R*

#### Stratification

#### Stratum graph

nodes IDB predicates

edges  $P \rightarrow Q$  if P depends on Q

label the edge with "-" if Q is a negated subgoal

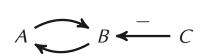
Stratifiable program: no cycle involving at least one negated edge

$$R(x) := S(x), \neg R(x)$$

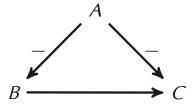


18/23

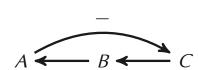
# More examples



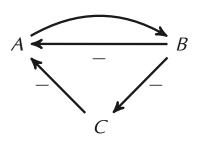
Stratifiable: A = B < C



Stratifiable:  $C \le B < A$ 



Not stratifiable:  $A \le B \le C < A$ 



Not stratifiable:  $A < B \le A$ 

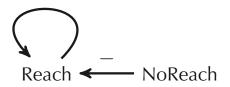
## Stratified example

Which target nodes cannot be reached from any source node?

 $NoReach(x) := Target(x), \neg Reach(x)$  (rule1)

Reach(x) := Source(x) (rule2)

Reach(x) := Reach(y), Link(y, x) (rule3)



Stratum 0 Source, Link, Target

Stratum 1 Reach

Stratum 2 NoReach

20/23

# Evaluation of stratified programs

*P* partitioned into a sequence  $P_1, \ldots, P_n$ 

Gives us an order in which to apply (each group of) rules

At each iteration k, execute each subprogram in sequence

- (1) Apply all the rules in  $P_1$
- (2) Apply all the rules in  $P_2$

:

(*n*) Apply all the rules in  $P_n$ 

# Evaluation of stratified programs

$$Reach(x) := Source(x)$$
 (P<sub>1</sub>)

$$Reach(x) := Reach(y), Link(y, x)$$
 (P<sub>1</sub>)

$$NoReach(x) := Target(x), \neg Reach(x)$$
 (P<sub>2</sub>)

**EDB** 

Source	Li	nk	Target
1	1	2	2
2	3	4	3
	2	4	4

**IDB** 

Reach	NoReach	
1	3	
2		
4		

22/23

# Further remarks on SQL recursion

- Requires stratified negation
- Only linear recursion

#### **Problems**

Arithmetic operations

introduce new values not present in the database

Multiset semantics

rules must be applied several times cycles in the data must be detected