

Datalog

Thursday, April 18, 2019 5:52 PM

- **Datalog**

- Another query language for relational model
- extends relational queries with recursion.

e.g. Souffle

- **Souffle**

- "splat" of Datalog
- set-based.

- **Facts**: tuples in database

Rules: queries

- data type { number
 symbol }

. decl **table**(name:type) create tables

table(data). insert table.

e.g. $Q1(y) :- \text{Movies}(x, y, z), z = 1940.$

order of variables matters

\Leftrightarrow name doesn't matter, unless you want to use later.

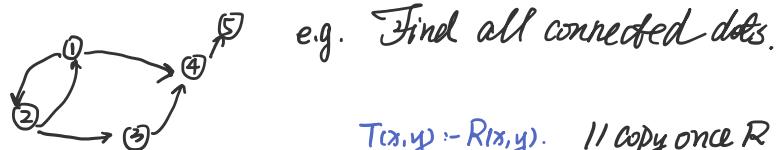
select $\underline{Q1(y)} := \text{Movies}(_, y, _), z=1940.$
 return y underscore means it's unnecessary.
 head (subgoal)
 atom "AND"
 $\underline{Q2(f, l)} := \text{Actor}(z, f, l), \text{Casts}(z, x), \text{Movie}(x, _, 1940).$
 head variables body
 ", " \Leftrightarrow "AND" (in some DBMS)

- **Extensional Database Predicates (EDB)** Predicates that are stored in databases.

- **Intentional Database Predicates (IDB)** Predicates that are expressed by rules.

- Multiple rules for same IDB means OR.

* e.g. Use graph to express friends. Find all connected friends.



$T(x, y) := R(x, y).$ // copy once R

$\underline{T(x, y) := R(x, y), T(z, y)}.$ // join $R \& T$.
 ↓
 Stop until the new $T(x, y)$ has covered all paths.

- **Datalog Semantics**

- **Fixpoint Semantic**

Start

$IDB_0 = \text{empty relations}$

$t=0.$

↓ Recur

11pm

$$IDB_{t+1} = \text{Compute Rules}(EDB, IDB_t)$$

$$t = t + 1$$

$$\text{Until } IDB_t = IDB_{t-1}$$

- Notice that since rules are monotone, $\phi = IDB_0 \subseteq IDB_1 \subseteq IDB_2 \subseteq \dots$

- $\begin{cases} \text{Right linear} & T(x,y) :- R(x,z), T(z,y). \\ \text{Left linear} & T(x,y) :- T(x,y), R(x,z). \\ \text{Non linear} & T(x,y) :- T(x,z), T(z,y). // \text{will terminate because set semantics.} \end{cases}$
- Non linear have least iterations.

- Model of data program

- A relation instance T is called a **model** if it satisfies these

$$\forall x, y \quad R(x,y) \rightarrow T(x,y).$$

$$T(z,y) :- R(x,y).$$

$$\forall x, y, z, \quad (R(x,z) \wedge T(z,y)) \rightarrow T(x,y). \quad T(x,y) :- R(x,z), T(z,y).$$

$$\Leftrightarrow \forall x, y, \exists z \text{ s.t. } R(x,z) \wedge T(z,y) \rightarrow T(x,y)$$

The datalog program computes minimum model.

o Aggregates aggregate-name variable : { relation to compute aggregate on }

count, min, max, sum

e.g. $\Delta(c) :- c = \text{count} : \{ \text{Actor}(-, y, -), y = 'John' \}$

count will ignore empty group.

e.g. $\Delta(\text{minId}) :- \underline{\text{minId}} = \text{min } x : \{ \text{Actor}(x, y, -), y = 'John' \}$
assign variable to the
value of the aggregate.

o Grouping

e.g. $\Delta(y, c) :- \text{Movie}(-, -, y), c = \text{count} : \{ \text{Movie}(-, -, y) \}$

o Negation "!"

e.g., find all descendants of Bob that are not of Alice.

$D(x, y) :- \text{ParentChild}(x, y)$.

$D(x, y) :- D(x, y), \text{ParentChild}(y, z)$. // calculate all descendants.

$\Delta(x) :- D("Bob", x), !D("Alice", x)$.

e.g. Compute children in the same generation. (two people are in the same gene if they are
descendants are at the same gen of common ancestor.)

$SG(x, y) :- \text{ParentChild}(p, x), \text{ParentChild}(p, y), x < y$ // common parent

$SG(x, y) :- \text{ParentChild}(p, x), \text{ParentChild}(q, y), SG(p, q), \underline{x < y}$ // parents at same generation
avoid output duplicate

- A datalog rule is **safe** if every variable appears in some **positive, non-aggregated relational atom**.

e.g. $\text{UI}(\text{x}, \text{y}) :- \text{ParentChild}(\text{"Alice"}, \text{x}), \text{y} \neq \text{"Bob"} // \text{infinite} :: \text{did not specify y}.$

`U2(x,y) :- ParentChild("Alice", x), !ParentChild(x,y) // did not specify y`

because!
so y can be any string \Rightarrow not safe!

- Stratified Datalog

- 0 Recursion does not cope well with aggregates or negation.

e.g., $A() :- !, B();$
 $B() :- !, A();$

- o A datalog program is **stratified** if it can be partitioned into strata.

Sin. *Stratum*

- Only IDB predicates defined in strata $1, 2, \dots, n$ may appear under !

or age in stratum 1.

(Many datalog DBMS accept only sorted data)

e.g., $D(x,y) :- \text{ParentChild}(x,y).$ // Stratum I
 $D(x,y) :- \text{ParentChild}(x,z), D(z,y).$

```
Q(X) :- D("Alice", X), !("Bob", X). // Stratum II
```

- If we don't use agg or negation, then the Datalog program is already stratified.

