

T1: Data models and query languages

L7: Datalog evaluation

Wolfgang Gatterbauer

CS7240 Principles of scalable data management (sp20)

<https://northeastern-datalab.github.io/cs7240/sp20/>

Version 1/28/2020

The issue of associativity

null appears in a join column. No matter what choice is taken, \bowtie is not associative. Consider the relations

q	r	s
$\begin{array}{ c c } \hline A & B \\ \hline 1 & 2 \\ \hline \end{array}$	$\begin{array}{ c c } \hline B & C \\ \hline 2 & 3 \\ \hline \end{array}$	$\begin{array}{ c c } \hline A & C \\ \hline 1 & 4 \\ \hline \end{array}$

Computing $(q \bowtie r) \bowtie s$ we get

q'
$\begin{array}{ c c c } \hline A & B & C \\ \hline 1 & 2 & 3 \\ \hline 1 & \perp & 4 \\ \hline \end{array}$

while $q \bowtie (r \bowtie s)$ gives

q''
$\begin{array}{ c c c } \hline A & B & C \\ \hline 1 & 2 & 4 \\ \hline \perp & 2 & 3 \\ \hline \end{array}$

$$\{a^-\} \bar{\cup}_a (\{a\} \bar{\cup}_a \{b\}) = \{a^-\}$$

$$(\{a^-\} \bar{\cup}_a \{a\}) \bar{\cup}_a \{b\} = \{b\}$$

left outer join example from p392 in "Maier. The theory of relational databases, 1983."

right preferred union example from "Gatterbauer, Suciu. Conflict resolution using trust mapping. SIGMOD 2010."

Outline: Datalog

- Datalog
 - Datalog rules
 - Recursion
 - Semantics
 - Datalog[¬]: Negation, stratification
 - Datalog[±]
 - Stable model semantics (Answer set programming)
 - Datalog vs. RA
 - Naive and Semi-naive evaluation

RA to Datalog by examples: Union



R(A,B,C)
S(D,E,F)
T(G,H)

RA:

$R(A,B,C) \cup S(D,E,F)$

Datalog:

?

RA to Datalog by examples: Union



R(A,B,C)
S(D,E,F)
T(G,H)

RA:

$$R(A,B,C) \cup S(D,E,F)$$

Datalog:

$$U(x,y,z) :- R(x,y,z)$$
$$U(x,y,z) :- S(x,y,z)$$

RA to Datalog by examples: Intersection



R(A,B,C)
S(D,E,F)
T(G,H)

RA:

$$R(A,B,C) \cap S(D,E,F)$$

Datalog:

?

RA to Datalog by examples: Intersection



R(A,B,C)
S(D,E,F)
T(G,H)

RA:

$$R(A,B,C) \cap S(D,E,F)$$

Datalog:

$$I(x,y,z) :- R(x,y,z), S(x,y,z)$$

RA to Datalog by examples: Selection



R(A,B,C)
S(D,E,F)
T(G,H)

RA:

$\sigma_{x>100 \text{ and } y='foo'} (R)$

Datalog:

?

RA to Datalog by examples: Selection



R(A,B,C)
S(D,E,F)
T(G,H)

RA:

$\sigma_{x>100 \text{ and } y='foo'} (R)$

Datalog:

$S(x,y,z) :- R(x,y,z), x > 100, y='foo'$

Handwritten annotations: A red arrow points from the 'y' in the Datalog query to the 'y' in the RA query. A red circle highlights 'y='foo'' in the Datalog query. A red checkmark is above the arrow, and a red 'X' is above the circle.

$\sigma_{x>100 \text{ or } y='foo'} (R)$

?

RA to Datalog by examples: Selection



R(A,B,C)
S(D,E,F)
T(G,H)

RA:

$$\sigma_{x>100 \text{ and } y='foo'} (R)$$

Datalog:

$$S(x,y,z) :- R(x,y,z), x > 100, y='foo'$$
$$\sigma_{x>100 \text{ or } y='foo'} (R)$$
$$S(x,y,z) :- R(x,y,z), x > 100$$
$$S(x,y,z) :- R(x,y,z), y='foo'$$

RA to Datalog by examples: Equi-join



R(A,B,C)
S(D,E,F)
T(G,H)

RA:

$R \bowtie_{R.A=S.D \text{ and } R.B=S.E} S$

Datalog:

?

RA to Datalog by examples: Equi-join



R(A,B,C)
S(D,E,F)
T(G,H)

RA:

$$R \bowtie_{R.A=S.D \text{ and } R.B=S.E} S$$

Datalog:

$$J(x,y,z,q) :- R(x,y,z), S(x,y,q)$$

RA to Datalog by examples: Projection



R(A,B,C)
S(D,E,F)
T(G,H)

RA:

$$\pi_A(R)$$

Datalog:

?

RA to Datalog by examples: Projection



R(A,B,C)
S(D,E,F)
T(G,H)

EDB →

RA:

$\pi_A(R)$

Datalog:

$P(x) :- R(x, \overline{y}, \overline{z})$

IPB

RA to Datalog by examples: Difference



R(A,B,C)
S(D,E,F)
T(G,H)

RA:

R-S

Datalog:

?

RA to Datalog by examples: Difference



R(A,B,C)
S(D,E,F)
T(G,H)

RA:

R-S



Datalog⁻: (we need to add **negation**)

$D(x,y,z) :- \text{R}(x,y,z), \text{not } S(x,y,z)$

Outline: Datalog

- Datalog
 - Datalog rules
 - Recursion
 - Semantics
 - Datalog[¬]: Negation, stratification
 - Datalog[±]
 - Stable model semantics (Answer set programming)
 - Datalog vs. RA
 - Naive and Semi-naive evaluation

Datalog Evaluation Algorithms

- Needs to preserve the efficiency of query optimizers, while extending them to recursion
- Two general strategies:
 - Naive datalog evaluation
 - Semi-naive datalog evaluation
- Some more powerful optimizations:
 - Magic sets (not covered)

Naive Datalog evaluation

$\text{Path}(x,y) \text{ :- Edge}(x,y)$
 $\text{Path}(x,z) \text{ :- Edge}(x,y), \text{Path}(y,z)$

$\text{Path}^{(0)} := \emptyset, t:=0$

Repeat {

 inc(t)

$\text{Path}^{(t)}(x,y) := \text{Edge}(x,y) \cup \Pi_{xy}(\text{Edge}(x,z) \bowtie \text{Path}^{(t-1)}(z,y))$

"immediate consequence operator"

until $\text{Path}^{(t)} = \text{Path}^{(t-1)}$ }

- The same facts are discovered over and over again
- Next: The **semi-naive** algorithm tries to reduce the number of facts discovered multiple times

Background: Incremental View Maintenance

Let V be a view computed by one datalog rule (no recursion)

View :- body

If (some of) the relations are updated:

$$R_1 \leftarrow R_1 \cup \Delta R_1, R_2 \leftarrow R_2 \cup \Delta R_2, \dots$$

Then the view is also modified as follows:

$$V \leftarrow V \cup \Delta V$$

Incremental view maintenance:

Compute ΔV without having to recompute V

Background: Incremental View Maintenance



Example 1:

$V(x,y) :- R(x,z), S(z,y)$

If $R \leftarrow RU\Delta R$,
then what is $\Delta V(x,y)$?

?

Background: Incremental View Maintenance



Example 1:

$$V(x,y) :- R(x,z), S(z,y)$$

$$\Delta V(x,y) :- \Delta R(x,z), S(z,y)$$

If $R \leftarrow RU\Delta R$,
then what is $\Delta V(x,y)$?

Algebra:

$$V = R \bowtie S$$

$$VU\Delta V = (RU\Delta R) \bowtie S$$

?

Background: Incremental View Maintenance



Example 1:

$$V(x,y) :- R(x,z), S(z,y)$$

$$\Delta V(x,y) :- \Delta R(x,z), S(z,y)$$

If $R \leftarrow RU\Delta R$,
then what is $\Delta V(x,y)$?

$$z = x \cdot y$$

$$z + \Delta z = (x + \Delta x) \cdot y$$

?

Algebra:

$$V = R \bowtie S$$

$$V \cup \Delta V = (R \cup \Delta R) \bowtie S$$

$$V \cup \Delta V = (R \bowtie S) \cup (\Delta R \bowtie S)$$

$$V \cup \Delta V = V \cup (\Delta R \bowtie S)$$

$$\Delta V = \Delta R \bowtie S$$

*Union
distributes
over join*

Background: Incremental View Maintenance



Example 1:

$$V(x,y) :- R(x,z), S(z,y)$$

$$\Delta V(x,y) :- \Delta R(x,z), S(z,y)$$

If $R \leftarrow RU\Delta R$,
then what is $\Delta V(x,y)$?

\otimes distributes over \oplus

$$z = x \cdot y$$

$$z + \Delta z = (x + \Delta x) \cdot y$$

...

$$\Delta z = \Delta x \cdot y$$

Algebra:

$$V = R \bowtie S$$

$$V \cup \Delta V = (R \cup \Delta R) \bowtie S$$

$$V \cup \Delta V = (R \bowtie S) \cup (\Delta R \bowtie S)$$

$$V \cup \Delta V = V \cup (\Delta R \bowtie S)$$

$$\Delta V = \Delta R \bowtie S$$

*Union
distributes
over join*

Background: Incremental View Maintenance



Example 2:

$V(x,y) :- R(x,z), S(z,y)$

If $R \leftarrow RU\Delta R$, and $S \leftarrow SU\Delta S$,
then what is $\Delta V(x,y)$?

?

Background: Incremental View Maintenance



Example 2:

$$V(x,y) :- R(x,z), S(z,y)$$

$$\Delta V(x,y) :- \Delta R(x,z), S(z,y)$$

$$\Delta V(x,y) :- R(x,z), \Delta S(z,y)$$

$$\Delta V(x,y) :- \Delta R(x,z), \Delta S(z,y)$$

If $R \leftarrow RU\Delta R$, and $S \leftarrow SU\Delta S$,
then what is $\Delta V(x,y)$?

\otimes distributes over \oplus

$$z = x \cdot y$$

$$z + \Delta z = (x + \Delta x) \cdot (y + \Delta y)$$

?

Algebra:

$$V = R \bowtie S$$

$$V \cup \Delta V = (RU\Delta R) \bowtie (SU\Delta S)$$

?

Union
distributes
over join

Background: Incremental View Maintenance



Example 2:

$$V(x,y) :- R(x,z), S(z,y)$$

If $R \leftarrow RU\Delta R$, and $S \leftarrow SU\Delta S$,
then what is $\Delta V(x,y)$?

$$\Delta V(x,y) :- \Delta R(x,z), S(z,y)$$

$$\Delta V(x,y) :- R(x,z), \Delta S(z,y)$$

$$\Delta V(x,y) :- \Delta R(x,z), \Delta S(z,y)$$

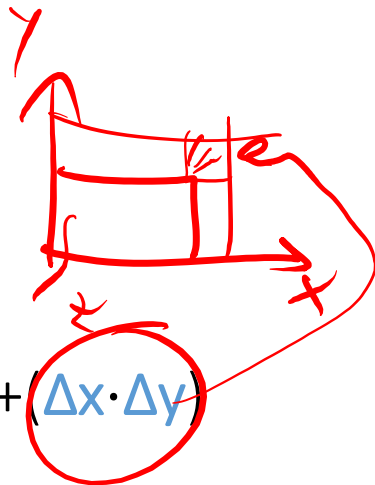
\otimes distributes over \oplus

$$z = x \cdot y$$

$$z + \Delta z = (x + \Delta x) \cdot (y + \Delta y)$$

...

$$\Delta z = (\Delta x \cdot y) + (x \cdot \Delta y) + (\Delta x \cdot \Delta y)$$



Algebra:

$$V = R \bowtie S$$

$$V \cup \Delta V = (RU\Delta R) \bowtie (SU\Delta S)$$

$$V \cup \Delta V = (R \bowtie S) \cup (\Delta R \bowtie S) \cup (R \bowtie \Delta S) \cup (\Delta R \bowtie \Delta S)$$

$$V \cup \Delta V = V \cup (\Delta R \bowtie S) \cup (R \bowtie \Delta S) \cup (\Delta R \bowtie \Delta S)$$

$$\Delta V = (\Delta R \bowtie S) \cup (R \bowtie \Delta S) \cup (\Delta R \bowtie \Delta S)$$

Union
distributes
over join

Background: Incremental View Maintenance



Example 3:

$V(x,y) :- R(x,z), R(z,y)$

If $R \leftarrow RU\Delta R$,
then what is $\Delta V(x,y)$?

?

Background: Incremental View Maintenance



Example 3:

$$V(x,y) :- R(x,z), R(z,y)$$

$$\Delta V(x,y) :- \Delta R(x,z), R(z,y)$$

$$\Delta V(x,y) :- R(x,z), \Delta R(z,y)$$

$$\Delta V(x,y) :- \Delta R(x,z), \Delta R(z,y)$$

If $R \leftarrow RU\Delta R$,
then what is $\Delta V(x,y)$?

\otimes distribute over \oplus

$$z = x^2$$

$$z + \Delta z = (x + \Delta x)^2$$

...

$$\Delta z = 2x\Delta x + \Delta x^2$$

Algebra:

$$V = R \bowtie R$$

$$V \cup \Delta V = (RU\Delta R) \bowtie (SU\Delta R)$$

$$V \cup \Delta V = (R \bowtie R) \cup (\Delta R \bowtie R) \cup (R \bowtie \Delta R) \cup (\Delta R \bowtie \Delta R)$$

$$V \cup \Delta V = V \cup (\Delta R \bowtie R) \cup (R \bowtie \Delta R) \cup (\Delta R \bowtie \Delta R)$$

$$\Delta V = (\Delta R \bowtie R) \cup (R \bowtie \Delta R) \cup (\Delta R \bowtie \Delta R)$$

Union
distributes
over join

Semi-Naive Datalog evaluation

```
Path(x,y) :- Edge(x,y)
Path(x,z) :- Edge(x,y), Path(y,z)
```

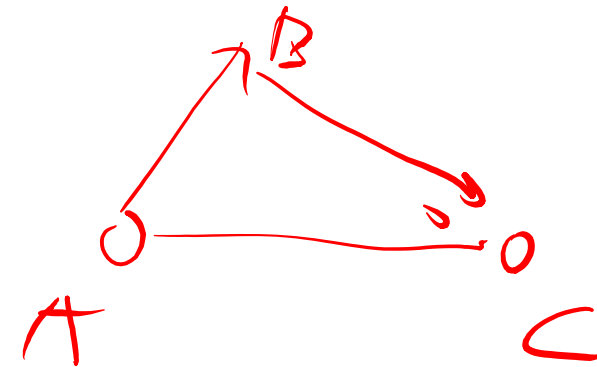
```
Path := Edge(x,z); ΔPath := Edge(x,z)
```

```
Repeat {
```

```
    ΔPath(x,y) :=  $\Pi_{xy}(\text{Edge}(x,z) \bowtie \Delta\text{Path}(z,y)) - \text{Path}(x,y)$ 
```

```
    Path := Path  $\cup$  ΔPath
```

```
until ΔPath =  $\emptyset$ }
```



Datalog Summary

- EDB (base relations) and IDB (derived relations)
- Datalog program = set of rules
- Datalog is recursive

- Some reminders about semantics:
 - Multiple atoms in a rule mean join (or intersection)
 - Variables with the same name are join variables
 - Multiple rules with same head mean union

Datalog and SQL

- Stratified data (w/ recursion, w/o +, *, ...): expresses precisely* queries in PTIME
 - Cannot find a Hamiltonian cycle (why?)
- SQL has also been extended to express recursive queries:
 - Use a recursive "with" clause, also CTE (Common Table Expression)
 - Often with bizarre restrictions...
 - ... Just use datalog

* needs to use the < predicate