

DATABASE THEORY

Lecture 14: Datalog Evaluation

Markus Krötzsch

Knowledge-Based Systems

TU Dresden, 16 June 2025

More recent versions of this slide deck might be available.
For the most current version of this course, see
https://iccl.inf.tu-dresden.de/web/Database_Theory/en

Review: Datalog

A rule-based recursive query language

```
father(alice, bob)
```

```
mother(alice, carla)
```

```
Parent(x, y) ← father(x, y)
```

```
Parent(x, y) ← mother(x, y)
```

```
SameGeneration(x, x)
```

```
SameGeneration(x, y) ← Parent(x, v) ∧ Parent(y, w) ∧ SameGeneration(v, w)
```

- Datalog is more complex than FO query answering
- Datalog is more expressive than FO query answering
- Semipositive Datalog with a successor ordering captures P
- Datalog containment is undecidable

Remaining question: **How can Datalog query answering be implemented?**

Implementing Datalog

FO queries (and thus also CQs and UCQs) are supported by almost all DBMS
~ many specific implementation and optimisation techniques

How can Datalog queries be answered in practice?

~ techniques for dealing with recursion in DBMS query answering

Implementing Datalog

FO queries (and thus also CQs and UCQs) are supported by almost all DBMS
~ many specific implementation and optimisation techniques

How can Datalog queries be answered in practice?

~ techniques for dealing with recursion in DBMS query answering

There are two major paradigms for answering recursive queries:

- **Bottom-up:** derive conclusions by applying rules to given facts
- **Top-down:** search for proofs to infer results given query

Computing Datalog Query Answers Bottom-Up

We already saw a way to compute Datalog answers bottom-up:
the step-wise computation of the consequence operator T_P

Bottom-up computation is known under many names:

- **Forward-chaining** since rules are “chained” from premise to conclusion
(common in logic programming)
- **Materialisation** since inferred facts are stored (“materialised”)
(common in databases)
- **Saturation** since the input database is “saturated” with inferences
(common in theorem proving)
- **Deductive closure** since we “close” the input under entailments
(common in formal logic)

Naive Evaluation of Datalog Queries

A direct approach for computing T_P^∞

```
01   $T_P^0 := \emptyset$ 
02   $i := 0$ 
03  repeat :
04       $T_P^{i+1} := \emptyset$ 
05      for  $H \leftarrow B_1 \wedge \dots \wedge B_\ell \in P :$ 
06          for  $\theta \in B_1 \wedge \dots \wedge B_\ell(T_P^i) :$ 
07               $T_P^{i+1} := T_P^{i+1} \cup \{H\theta\}$ 
08       $i := i + 1$ 
09  until  $T_P^{i-1} = T_P^i$ 
10  return  $T_P^i$ 
```

Notation for line 06/07:

- a substitution θ is a mapping from variables to database elements
- for a formula F , we write $F\theta$ for the formula obtained by replacing each free variable x in F by $\theta(x)$
- for a CQ Q and database \mathcal{I} , we write $\theta \in Q(\mathcal{I})$ if $\mathcal{I} \models Q\theta$

What's Wrong with Naive Evaluation?

An example Datalog program:

```
e(1, 2)  e(2, 3)  e(3, 4)  e(4, 5)
(R1)      T(x, y) ← e(x, y)
(R2)      T(x, z) ← T(x, y) ∧ T(y, z)
```

$$T_P^0 = \emptyset$$

$$T_P^1 = \{T(1, 2), T(2, 3), T(3, 4), T(4, 5)\}$$

$$T_P^2 = T_P^1 \cup \{T(1, 3), T(2, 4), T(3, 5)\}$$

$$T_P^3 = T_P^2 \cup \{T(1, 4), T(2, 5), T(1, 5)\}$$

$$T_P^4 = T_P^3 = T_P^\infty$$

What's Wrong with Naive Evaluation?

An example Datalog program:

| | | | | |
|------|---|---------|---------|---------|
| | e(1, 2) | e(2, 3) | e(3, 4) | e(4, 5) |
| (R1) | T(x, y) \leftarrow e(x, y) | | | |
| (R2) | T(x, z) \leftarrow T(x, y) \wedge T(y, z) | | | |

How many body matches do we need to iterate over?

$T_P^0 = \emptyset$ initialisation

$T_P^1 = \{T(1, 2), T(2, 3), T(3, 4), T(4, 5)\}$

$T_P^2 = T_P^1 \cup \{T(1, 3), T(2, 4), T(3, 5)\}$

$T_P^3 = T_P^2 \cup \{T(1, 4), T(2, 5), T(1, 5)\}$

$T_P^4 = T_P^3 = T_P^\infty$

What's Wrong with Naive Evaluation?

An example Datalog program:

| | | | | |
|------|---|---------|---------|---------|
| | e(1, 2) | e(2, 3) | e(3, 4) | e(4, 5) |
| (R1) | T(x, y) \leftarrow e(x, y) | | | |
| (R2) | T(x, z) \leftarrow T(x, y) \wedge T(y, z) | | | |

How many body matches do we need to iterate over?

$T_P^0 = \emptyset$ initialisation

$T_P^1 = \{T(1, 2), T(2, 3), T(3, 4), T(4, 5)\}$ 4 matches for (R1)

$T_P^2 = T_P^1 \cup \{T(1, 3), T(2, 4), T(3, 5)\}$

$T_P^3 = T_P^2 \cup \{T(1, 4), T(2, 5), T(1, 5)\}$

$T_P^4 = T_P^3 = T_P^\infty$

What's Wrong with Naive Evaluation?

An example Datalog program:

| | | | | |
|------|---|---------|---------|---------|
| | e(1, 2) | e(2, 3) | e(3, 4) | e(4, 5) |
| (R1) | T(x, y) \leftarrow e(x, y) | | | |
| (R2) | T(x, z) \leftarrow T(x, y) \wedge T(y, z) | | | |

How many body matches do we need to iterate over?

$T_P^0 = \emptyset$ initialisation

$T_P^1 = \{T(1, 2), T(2, 3), T(3, 4), T(4, 5)\}$ 4 matches for (R1)

$T_P^2 = T_P^1 \cup \{T(1, 3), T(2, 4), T(3, 5)\}$ $4 \times (R1) + 3 \times (R2)$

$T_P^3 = T_P^2 \cup \{T(1, 4), T(2, 5), T(1, 5)\}$

$T_P^4 = T_P^3 = T_P^\infty$

What's Wrong with Naive Evaluation?

An example Datalog program:

| | | | | |
|------|---|---------|---------|---------|
| | e(1, 2) | e(2, 3) | e(3, 4) | e(4, 5) |
| (R1) | T(x, y) \leftarrow e(x, y) | | | |
| (R2) | T(x, z) \leftarrow T(x, y) \wedge T(y, z) | | | |

How many body matches do we need to iterate over?

$T_P^0 = \emptyset$ initialisation

$T_P^1 = \{T(1, 2), T(2, 3), T(3, 4), T(4, 5)\}$ 4 matches for (R1)

$T_P^2 = T_P^1 \cup \{T(1, 3), T(2, 4), T(3, 5)\}$ $4 \times (R1) + 3 \times (R2)$

$T_P^3 = T_P^2 \cup \{T(1, 4), T(2, 5), T(1, 5)\}$ $4 \times (R1) + 8 \times (R2)$

$T_P^4 = T_P^3 = T_P^\infty$

What's Wrong with Naive Evaluation?

An example Datalog program:

| | | | | |
|------|---|---------|---------|---------|
| | e(1, 2) | e(2, 3) | e(3, 4) | e(4, 5) |
| (R1) | T(x, y) \leftarrow e(x, y) | | | |
| (R2) | T(x, z) \leftarrow T(x, y) \wedge T(y, z) | | | |

How many body matches do we need to iterate over?

$T_P^0 = \emptyset$ initialisation

$T_P^1 = \{T(1, 2), T(2, 3), T(3, 4), T(4, 5)\}$ 4 matches for (R1)

$T_P^2 = T_P^1 \cup \{T(1, 3), T(2, 4), T(3, 5)\}$ $4 \times (R1) + 3 \times (R2)$

$T_P^3 = T_P^2 \cup \{T(1, 4), T(2, 5), T(1, 5)\}$ $4 \times (R1) + 8 \times (R2)$

$T_P^4 = T_P^3 = T_P^\infty$ $4 \times (R1) + 10 \times (R2)$

What's Wrong with Naive Evaluation?

An example Datalog program:

| | | | | |
|------|---|---------|---------|---------|
| | e(1, 2) | e(2, 3) | e(3, 4) | e(4, 5) |
| (R1) | T(x, y) \leftarrow e(x, y) | | | |
| (R2) | T(x, z) \leftarrow T(x, y) \wedge T(y, z) | | | |

How many body matches do we need to iterate over?

$T_P^0 = \emptyset$ initialisation

$T_P^1 = \{T(1, 2), T(2, 3), T(3, 4), T(4, 5)\}$ 4 matches for (R1)

$T_P^2 = T_P^1 \cup \{T(1, 3), T(2, 4), T(3, 5)\}$ $4 \times (R1) + 3 \times (R2)$

$T_P^3 = T_P^2 \cup \{T(1, 4), T(2, 5), T(1, 5)\}$ $4 \times (R1) + 8 \times (R2)$

$T_P^4 = T_P^3 = T_P^\infty$ $4 \times (R1) + 10 \times (R2)$

In total, we considered 37 matches to derive 11 facts

Less Naive Evaluation Strategies

Does it really matter how often we consider a rule match?

After all, each fact is added only once ...

Less Naive Evaluation Strategies

Does it really matter how often we consider a rule match?

After all, each fact is added only once ...

In practice, finding applicable rules takes significant time, even if the conclusion does not need to be added – iteration takes time!

~ huge potential for optimisation

Observation:

we derive the same conclusions over and over again in each step

Idea: apply rules only to newly derived facts

~ semi-naive evaluation

Semi-Naive Evaluation

The computation yields sets $T_P^0 \subseteq T_P^1 \subseteq T_P^2 \subseteq \dots \subseteq T_P^\infty$

- For an IDB predicate R , let R^i be the “predicate” that contains exactly the R -facts in T_P^i
- For $i \leq 1$, let Δ_R^i be the collection of facts $R^i \setminus R^{i-1}$

We can restrict rules to use only some computations.

Semi-Naive Evaluation

The computation yields sets $T_P^0 \subseteq T_P^1 \subseteq T_P^2 \subseteq \dots \subseteq T_P^\infty$

- For an IDB predicate R , let R^i be the “predicate” that contains exactly the R -facts in T_P^i
- For $i \leq 1$, let Δ_R^i be the collection of facts $R^i \setminus R^{i-1}$

We can restrict rules to use only some computations.

Some options for the computation in step $i + 1$:

$T(x, z) \leftarrow T^i(x, y) \wedge T^i(y, z)$ same as original rule

$T(x, z) \leftarrow \Delta_T^i(x, y) \wedge \Delta_T^i(y, z)$ restrict to new facts

$T(x, z) \leftarrow \Delta_T^i(x, y) \wedge T^i(y, z)$ partially restrict to new facts

$T(x, z) \leftarrow T^i(x, y) \wedge \Delta_T^i(y, z)$ partially restrict to new facts

What to choose?

Semi-Naive Evaluation (2)

Inferences that involve new and old facts are necessary:

$e(1, 2) \quad e(2, 3) \quad e(3, 4) \quad e(4, 5)$

(R1) $T(x, y) \leftarrow e(x, y)$

(R2) $T(x, z) \leftarrow T(x, y) \wedge T(y, z)$

$$T_P^0 = \emptyset$$

$$\Delta_T^1 = \{T(1, 2), T(2, 3), T(3, 4), T(4, 5)\} \quad T_P^1 = \Delta_T^1$$

$$\Delta_T^2 = \{T(1, 3), T(2, 4), T(3, 5)\} \quad T_P^2 = T_P^1 \cup \Delta_T^2$$

$$\Delta_T^3 = \{T(1, 4), T(2, 5), T(1, 5)\} \quad T_P^3 = T_P^2 \cup \Delta_T^3$$

$$\Delta_T^4 = \emptyset \quad T_P^4 = T_P^3 = T_P^\infty$$

Semi-Naive Evaluation (2)

Inferences that involve new and old facts are necessary:

$$\begin{array}{ll} e(1, 2) & e(2, 3) \quad e(3, 4) \quad e(4, 5) \\ (R1) \quad T(x, y) \leftarrow e(x, y) \\ (R2) \quad T(x, z) \leftarrow T(x, y) \wedge T(y, z) \end{array}$$

$$T_P^0 = \emptyset$$

$$\Delta_T^1 = \{T(1, 2), T(2, 3), T(3, 4), T(4, 5)\} \quad T_P^1 = \Delta_T^1$$

$$\Delta_T^2 = \{T(1, 3), T(2, 4), T(3, 5)\} \quad T_P^2 = T_P^1 \cup \Delta_T^2$$

$$\Delta_T^3 = \{T(1, 4), T(2, 5), T(1, 5)\} \quad T_P^3 = T_P^2 \cup \Delta_T^3$$

$$\Delta_T^4 = \emptyset \quad T_P^4 = T_P^3 = T_P^\infty$$

To derive $T(1, 4)$ in Δ_T^3 , we need to combine

$T(1, 3) \in \Delta_T^2$ with $T(3, 4) \in \Delta_T^1$ or $T(1, 2) \in \Delta_T^1$ with $T(2, 4) \in \Delta_T^2$

Semi-Naive Evaluation (2)

Inferences that involve new and old facts are necessary:

$$\begin{array}{ll} e(1, 2) & e(2, 3) \quad e(3, 4) \quad e(4, 5) \\ (R1) \quad T(x, y) \leftarrow e(x, y) \\ (R2) \quad T(x, z) \leftarrow T(x, y) \wedge T(y, z) \end{array}$$

$$T_P^0 = \emptyset$$

$$\Delta_T^1 = \{T(1, 2), T(2, 3), T(3, 4), T(4, 5)\} \quad T_P^1 = \Delta_T^1$$

$$\Delta_T^2 = \{T(1, 3), T(2, 4), T(3, 5)\} \quad T_P^2 = T_P^1 \cup \Delta_T^2$$

$$\Delta_T^3 = \{T(1, 4), T(2, 5), T(1, 5)\} \quad T_P^3 = T_P^2 \cup \Delta_T^3$$

$$\Delta_T^4 = \emptyset \quad T_P^4 = T_P^3 = T_P^\infty$$

To derive $T(1, 4)$ in Δ_T^3 , we need to combine

$T(1, 3) \in \Delta_T^2$ with $T(3, 4) \in \Delta_T^1$ or $T(1, 2) \in \Delta_T^1$ with $T(2, 4) \in \Delta_T^2$

↪ rule $T(x, z) \leftarrow \Delta_T^i(x, y) \wedge \Delta_T^i(y, z)$ is not enough

Semi-Naive Evaluation (3)

Correct approach: consider only rule application that use **at least one** newly derived IDB atom

For example program:

| | | | | |
|----------|------------------------------|--|--|-----------|
| | $e(1, 2)$ | $e(2, 3)$ | $e(3, 4)$ | $e(4, 5)$ |
| $(R1)$ | $T(x, y) \leftarrow e(x, y)$ | | | |
| $(R2.1)$ | | $T(x, z) \leftarrow \Delta_T^i(x, y) \wedge T^i(y, z)$ | | |
| $(R2.2)$ | | | $T(x, z) \leftarrow T^i(x, y) \wedge \Delta_T^i(y, z)$ | |

Semi-Naive Evaluation (3)

Correct approach: consider only rule application that use **at least one** newly derived IDB atom

For example program:

| | | | | |
|----------|------------------------------|--|--|-----------|
| | $e(1, 2)$ | $e(2, 3)$ | $e(3, 4)$ | $e(4, 5)$ |
| $(R1)$ | $T(x, y) \leftarrow e(x, y)$ | | | |
| $(R2.1)$ | | $T(x, z) \leftarrow \Delta_T^i(x, y) \wedge T^i(y, z)$ | | |
| $(R2.2)$ | | | $T(x, z) \leftarrow T^i(x, y) \wedge \Delta_T^i(y, z)$ | |

There is still redundancy here: the matches for $T(x, z) \leftarrow \Delta_T^i(x, y) \wedge \Delta_T^i(y, z)$ are covered by both $(R2.1)$ and $(R2.2)$

Semi-Naive Evaluation (3)

Correct approach: consider only rule application that use **at least one** newly derived IDB atom

For example program:

$$\begin{array}{ll} & \mathbf{e(1,2) \ e(2,3) \ e(3,4) \ e(4,5)} \\ (R1) & \mathbf{T(x,y) \leftarrow e(x,y)} \\ (R2.1) & \mathbf{T(x,z) \leftarrow \Delta_T^i(x,y) \wedge T^i(y,z)} \\ (R2.2) & \mathbf{T(x,z) \leftarrow T^i(x,y) \wedge \Delta_T^i(y,z)} \end{array}$$

There is still redundancy here: the matches for $\mathbf{T(x,z) \leftarrow \Delta_T^i(x,y) \wedge \Delta_T^i(y,z)}$ are covered by both (R2.1) and (R2.2)

~ replace (R2.2) by the following rule:

$$(R2.2') \quad \mathbf{T(x,z) \leftarrow T^{i-1}(x,y) \wedge \Delta_T^i(y,z)}$$

EDB atoms do not change, so their Δ would be \emptyset

~ ignore such rules after the first iteration

Semi-Naive Evaluation: Example

$e(1, 2) \quad e(2, 3) \quad e(3, 4) \quad e(4, 5)$

$(R1) \quad T(x, y) \leftarrow e(x, y)$

$(R2.1) \quad T(x, z) \leftarrow \Delta_T^i(x, y) \wedge T^i(y, z)$

$(R2.2') \quad T(x, z) \leftarrow T^{i-1}(x, y) \wedge \Delta_T^i(y, z)$

$$T_P^0 = \emptyset$$

$$T_P^1 = \{T(1, 2), T(2, 3), T(3, 4), T(4, 5)\}$$

$$T_P^2 = T_P^1 \cup \{T(1, 3), T(2, 4), T(3, 5)\}$$

$$T_P^3 = T_P^2 \cup \{T(1, 4), T(2, 5), T(1, 5)\}$$

$$T_P^4 = T_P^3 = T_P^\infty$$

Semi-Naive Evaluation: Example

| | | | | |
|---------|--|---------|---------|---------|
| | e(1, 2) | e(2, 3) | e(3, 4) | e(4, 5) |
| (R1) | T(x, y) ← e(x, y) | | | |
| (R2.1) | T(x, z) ← Δ _T ⁱ (x, y) ∧ T ⁱ (y, z) | | | |
| (R2.2') | T(x, z) ← T ⁱ⁻¹ (x, y) ∧ Δ _T ⁱ (y, z) | | | |

How many body matches do we need to iterate over?

$$T_P^0 = \emptyset \quad \text{initialisation}$$

$$T_P^1 = \{T(1, 2), T(2, 3), T(3, 4), T(4, 5)\}$$

$$T_P^2 = T_P^1 \cup \{T(1, 3), T(2, 4), T(3, 5)\}$$

$$T_P^3 = T_P^2 \cup \{T(1, 4), T(2, 5), T(1, 5)\}$$

$$T_P^4 = T_P^3 = T_P^\infty$$

Semi-Naive Evaluation: Example

$e(1, 2) \quad e(2, 3) \quad e(3, 4) \quad e(4, 5)$

$(R1) \quad T(x, y) \leftarrow e(x, y)$

$(R2.1) \quad T(x, z) \leftarrow \Delta_T^i(x, y) \wedge T^i(y, z)$

$(R2.2') \quad T(x, z) \leftarrow T^{i-1}(x, y) \wedge \Delta_T^i(y, z)$

How many body matches do we need to iterate over?

$T_P^0 = \emptyset$ initialisation

$T_P^1 = \{T(1, 2), T(2, 3), T(3, 4), T(4, 5)\} \quad 4 \times (R1)$

$T_P^2 = T_P^1 \cup \{T(1, 3), T(2, 4), T(3, 5)\}$

$T_P^3 = T_P^2 \cup \{T(1, 4), T(2, 5), T(1, 5)\}$

$T_P^4 = T_P^3 = T_P^\infty$

Semi-Naive Evaluation: Example

| | | | | |
|---------|--|---------|---------|---------|
| | e(1, 2) | e(2, 3) | e(3, 4) | e(4, 5) |
| (R1) | T(x, y) ← e(x, y) | | | |
| (R2.1) | T(x, z) ← Δ _T ⁱ (x, y) ∧ T ⁱ (y, z) | | | |
| (R2.2') | T(x, z) ← T ⁱ⁻¹ (x, y) ∧ Δ _T ⁱ (y, z) | | | |

How many body matches do we need to iterate over?

$$T_P^0 = \emptyset \quad \text{initialisation}$$

$$T_P^1 = \{T(1, 2), T(2, 3), T(3, 4), T(4, 5)\} \quad 4 \times (R1)$$

$$T_P^2 = T_P^1 \cup \{T(1, 3), T(2, 4), T(3, 5)\} \quad 3 \times (R2.1)$$

$$T_P^3 = T_P^2 \cup \{T(1, 4), T(2, 5), T(1, 5)\}$$

$$T_P^4 = T_P^3 = T_P^{\infty}$$

Semi-Naive Evaluation: Example

| | | | | |
|---------|--|---------|---------|---------|
| | e(1, 2) | e(2, 3) | e(3, 4) | e(4, 5) |
| (R1) | T(x, y) ← e(x, y) | | | |
| (R2.1) | T(x, z) ← Δ _T ⁱ (x, y) ∧ T ⁱ (y, z) | | | |
| (R2.2') | T(x, z) ← T ⁱ⁻¹ (x, y) ∧ Δ _T ⁱ (y, z) | | | |

How many body matches do we need to iterate over?

$$T_P^0 = \emptyset \quad \text{initialisation}$$

$$T_P^1 = \{T(1, 2), T(2, 3), T(3, 4), T(4, 5)\} \quad 4 \times (R1)$$

$$T_P^2 = T_P^1 \cup \{T(1, 3), T(2, 4), T(3, 5)\} \quad 3 \times (R2.1)$$

$$T_P^3 = T_P^2 \cup \{T(1, 4), T(2, 5), T(1, 5)\} \quad 3 \times (R2.1), 2 \times (R2.2')$$

$$T_P^4 = T_P^3 = T_P^{\infty}$$

Semi-Naive Evaluation: Example

| | | | | |
|---------|--|---------|---------|---------|
| | e(1, 2) | e(2, 3) | e(3, 4) | e(4, 5) |
| (R1) | T(x, y) ← e(x, y) | | | |
| (R2.1) | T(x, z) ← Δ _T ⁱ (x, y) ∧ T ⁱ (y, z) | | | |
| (R2.2') | T(x, z) ← T ⁱ⁻¹ (x, y) ∧ Δ _T ⁱ (y, z) | | | |

How many body matches do we need to iterate over?

$$T_P^0 = \emptyset \quad \text{initialisation}$$

$$T_P^1 = \{T(1, 2), T(2, 3), T(3, 4), T(4, 5)\} \quad 4 \times (R1)$$

$$T_P^2 = T_P^1 \cup \{T(1, 3), T(2, 4), T(3, 5)\} \quad 3 \times (R2.1)$$

$$T_P^3 = T_P^2 \cup \{T(1, 4), T(2, 5), T(1, 5)\} \quad 3 \times (R2.1), 2 \times (R2.2')$$

$$T_P^4 = T_P^3 = T_P^{\infty} \quad 1 \times (R2.1), 1 \times (R2.2')$$

Semi-Naive Evaluation: Example

$$\begin{array}{c} \mathbf{e(1,2) \ e(2,3) \ e(3,4) \ e(4,5)} \\ (R1) \quad \mathbf{T(x,y) \leftarrow e(x,y)} \\ (R2.1) \quad \mathbf{T(x,z) \leftarrow \Delta_T^i(x,y) \wedge T^i(y,z)} \\ (R2.2') \quad \mathbf{T(x,z) \leftarrow T^{i-1}(x,y) \wedge \Delta_T^i(y,z)} \end{array}$$

How many body matches do we need to iterate over?

$$T_P^0 = \emptyset \quad \text{initialisation}$$

$$T_P^1 = \{\mathbf{T(1,2), T(2,3), T(3,4), T(4,5)}\} \quad 4 \times (R1)$$

$$T_P^2 = T_P^1 \cup \{\mathbf{T(1,3), T(2,4), T(3,5)}\} \quad 3 \times (R2.1)$$

$$T_P^3 = T_P^2 \cup \{\mathbf{T(1,4), T(2,5), T(1,5)}\} \quad 3 \times (R2.1), 2 \times (R2.2')$$

$$T_P^4 = T_P^3 = T_P^{\infty} \quad 1 \times (R2.1), 1 \times (R2.2')$$

In total, we considered 14 matches to derive 11 facts

Semi-Naive Evaluation: Full Definition

In general, a rule of the form

$$H(\vec{x}) \leftarrow e_1(\vec{y}_1) \wedge \dots \wedge e_n(\vec{y}_n) \wedge l_1(\vec{z}_1) \wedge l_2(\vec{z}_2) \wedge \dots \wedge l_m(\vec{z}_m)$$

is transformed into m rules

$$H(\vec{x}) \leftarrow e_1(\vec{y}_1) \wedge \dots \wedge e_n(\vec{y}_n) \wedge \Delta_{l_1}^i(\vec{z}_1) \wedge l_2^i(\vec{z}_2) \wedge \dots \wedge l_m^i(\vec{z}_m)$$

$$H(\vec{x}) \leftarrow e_1(\vec{y}_1) \wedge \dots \wedge e_n(\vec{y}_n) \wedge l_1^{i-1}(\vec{z}_1) \wedge \Delta_{l_2}^i(\vec{z}_2) \wedge \dots \wedge l_m^i(\vec{z}_m)$$

...

$$H(\vec{x}) \leftarrow e_1(\vec{y}_1) \wedge \dots \wedge e_n(\vec{y}_n) \wedge l_1^{i-1}(\vec{z}_1) \wedge l_2^{i-1}(\vec{z}_2) \wedge \dots \wedge \Delta_{l_m}^i(\vec{z}_m)$$

Advantages and disadvantages:

- Huge improvement over naive evaluation
- Some redundant computations remain (see example)
- Some overhead for implementation (store level of entailments)

Goal-Directed Datalog Evaluation

Top-Down Evaluation

Idea: we may not need to compute all derivations to answer a particular query

Example 14.1:

$$\begin{array}{cccc} e(1, 2) & e(2, 3) & e(3, 4) & e(4, 5) \\ (R1) & T(x, y) \leftarrow e(x, y) \\ (R2) & T(x, z) \leftarrow T(x, y) \wedge T(y, z) \\ \text{Query}(z) & \leftarrow T(2, z) \end{array}$$

The answers to Query are the T-successors of 2.

However, bottom-up computation would also produce facts like $T(1, 4)$, which are neither directly nor indirectly relevant for computing the query result.

Assumption

Assumption: For all techniques presented in this lecture, we assume that the given Datalog program is safe.

- This is without loss of generality (as shown in exercise).
- One can avoid this by adding more cases to algorithms.

Query-Subquery (QSQ)

QSQ is a technique for organising top-down Datalog query evaluation

Main principles:

- Apply backward chaining/resolution: start with query, find rules that can derive query, evaluate body atoms of those rules (subqueries) recursively
- Evaluate intermediate results “set-at-a-time” (using relational algebra on tables)
- Evaluate queries in a “data-driven” way, where operations are applied only to newly computed intermediate results (similar to idea in semi-naive evaluation)
- “Push” variable bindings (constants) from heads (queries) into bodies (subqueries)
- “Pass” variable bindings (constants) “sideways” from one body atom to the next

Details can be realised in several ways.

Adornments

To guide evaluation, we distinguish **free** and **bound** parameters in a predicate.

Example 14.2: If we want to derive atom $T(2, z)$ from the rule $T(x, z) \leftarrow T(x, y) \wedge T(y, z)$, then x will be bound to 2, while z is free.

Adornments

To guide evaluation, we distinguish **free** and **bound** parameters in a predicate.

Example 14.2: If we want to derive atom $T(2, z)$ from the rule $T(x, z) \leftarrow T(x, y) \wedge T(y, z)$, then x will be bound to 2, while z is free.

We use **adornments** to denote the free/bound parameters in predicates.

Example 14.3:

$$T^{bf}(x, z) \leftarrow T^{bf}(x, y) \wedge T^{bf}(y, z)$$

- since x is bound in the head, it is also bound in the first atom
- any match for the first atom binds y , so y is bound when evaluating the second atom (in left-to-right evaluation)

Adornments: Examples

The adornment of the head of a rule determines the adornments of the body atoms:

$$\begin{aligned} R^{bbb}(x, y, z) &\leftarrow R^{bbf}(x, y, v) \wedge R^{bbb}(x, v, z) \\ R^{fbf}(x, y, z) &\leftarrow R^{fbf}(x, y, v) \wedge R^{bbf}(x, v, z) \end{aligned}$$

Adornments: Examples

The adornment of the head of a rule determines the adornments of the body atoms:

$$\begin{aligned} R^{bbb}(x, y, z) &\leftarrow R^{bbf}(x, y, v) \wedge R^{bbb}(x, v, z) \\ R^{fbf}(x, y, z) &\leftarrow R^{fbf}(x, y, v) \wedge R^{bbf}(x, v, z) \end{aligned}$$

The order of body predicates affects the adornment:

$$\begin{aligned} S^{fff}(x, y, z) &\leftarrow T^{ff}(x, v) \wedge T^{ff}(y, w) \wedge R^{bbf}(v, w, z) \\ S^{fff}(x, y, z) &\leftarrow R^{fff}(v, w, z) \wedge T^{fb}(x, v) \wedge T^{fb}(y, w) \end{aligned}$$

~ For optimisation, some orders might be better than others

Auxiliary Relations for QSQ

To control evaluation, we store intermediate results in auxiliary relations.

When we “call” a rule with a head where some variables are bound, we need to provide the bindings as input

- ~ for adorned relation R^α , we use an auxiliary relation input_R^α
- ~ arity of input_R^α = number of b in α

The result of calling a rule should be the “completed” input, with values for the unbound variables added

- ~ for adorned relation R^α , we use an auxiliary relation output_R^α
- ~ arity of output_R^α = arity of R (= length of α)

Auxiliary Relations for QSQ (2)

When evaluating body atoms from left to right, we use supplementary relations sup_i

- ~ \rightsquigarrow bindings required to evaluate rest of rule after the i th body atom
- ~the first set of bindings sup_0 comes from input_R^α
- ~the last set of bindings sup_n go to output_R^α

Auxiliary Relations for QSQ (2)

When evaluating body atoms from left to right, we use supplementary relations sup_i

~ \rightsquigarrow bindings required to evaluate rest of rule after the i th body atom

~ \rightsquigarrow the first set of bindings sup_0 comes from input_R^a

~ \rightsquigarrow the last set of bindings sup_n go to output_R^a

Example 14.4:

$$\begin{array}{c} T^{bf}(x, z) \leftarrow T^{bf}(x, y) \wedge T^{bf}(y, z) \\ \uparrow \quad \nwarrow \uparrow \quad \Rightarrow \\ \text{input}_T^{bf} \Rightarrow \text{sup}_0[x] \quad \text{sup}_1[x, y] \quad \text{sup}_2[x, z] \Rightarrow \text{output}_T^{bf} \end{array}$$

- $\text{sup}_0[x]$ is copied from $\text{input}_T^{bf}[x]$ (with some exceptions, see exercise)
- $\text{sup}_1[x, y]$ is obtained by joining tables $\text{sup}_0[x]$ and $\text{output}_T^{bf}[x, y]$
- $\text{sup}_2[x, z]$ is obtained by joining tables $\text{sup}_1[x, y]$ and $\text{output}_T^{bf}[y, z]$
- $\text{output}_T^{bf}[x, z]$ is copied from $\text{sup}_2[x, z]$

(we use "named" notation like $[x, y]$ to suggest what to join on; the relations are the same)

QSQ Evaluation

The set of all auxiliary relations is called a **QSQ template** (for the given set of adorned rules)

General evaluation:

- add new tuples to auxiliary relations until reaching a fixed point
- evaluation of a rule can proceed as sketched on previous slide
- in addition, whenever new tuples are added to a sup relation that feeds into an IDB atom, the input relation of this atom is extended to include all binding given by sup (may trigger subquery evaluation)

~ there are many strategies for implementing this general scheme

QSQ Evaluation

The set of all auxiliary relations is called a **QSQ template** (for the given set of adorned rules)

General evaluation:

- add new tuples to auxiliary relations until reaching a fixed point
- evaluation of a rule can proceed as sketched on previous slide
- in addition, whenever new tuples are added to a sup relation that feeds into an IDB atom, the input relation of this atom is extended to include all binding given by sup (may trigger subquery evaluation)

~ there are many strategies for implementing this general scheme

Notation:

- for an EDB atom A , we write A^T for table that consists of all matches for A in the database

Recursive QSQ

Recursive QSQ (QSQR) takes a “depth-first” approach to QSQ

Evaluation of single rule in QSQR:

Given: adorned rule r with head predicate R^α ; current values of all QSQ relations

- (1) Copy tuples input_R^α (that unify with rule head) to sup_0^r
- (2) For each body atom A_1, \dots, A_n , do:
 - If A_i is an EDB atom, compute sup_i^r as projection of $\text{sup}_{i-1}^r \bowtie A_i^T$
 - If A_i is an IDB atom with adorned predicate S^β :
 - (a) Add new bindings from sup_{i-1}^r , combined with constants in A_i , to input_S^β
 - (b) If input_S^β changed, recursively evaluate all rules with head predicate S^β
 - (c) Compute sup_i^r as projection of $\text{sup}_{i-1}^r \bowtie \text{output}_S^\beta$
- (3) Add tuples in sup_n^r to output_R^α

Evaluation of query in QSQR:

Given: a Datalog program P and a conjunctive query $q[\vec{x}]$ (possibly with constants)

(1) Create an adorned program P^a :

- Turn the query $q[\vec{x}]$ into an adorned rule $\text{Query}^{ff\dots f}(\vec{x}) \leftarrow q[\vec{x}]$
- Recursively create adorned rules from rules in P for all adorned predicates in P^a .

(2) Initialise all auxiliary relations to empty sets.

(3) Evaluate the rule $\text{Query}^{ff\dots f}(\vec{x}) \leftarrow q[\vec{x}]$.

Repeat until no new tuples are added to any QSQ relation.

(4) Return output $\text{Query}^{ff\dots f}$.

QSQR Transformation: Example

Predicates S (same generation), p (parent), h (human)

$$S(x, x) \leftarrow h(x)$$

$$S(x, y) \leftarrow p(x, w) \wedge S(w, y) \wedge p(y, y)$$

with query $S(1, x)$.

QSQR Transformation: Example

Predicates S (same generation), p (parent), h (human)

$$S(x, x) \leftarrow h(x)$$

$$S(x, y) \leftarrow p(x, w) \wedge S(w, y) \wedge p(y, w)$$

with query $S(1, x)$.

~> Query rule: $\text{Query}(x) \leftarrow S(1, x)$

Transformed rules:

QSQR Transformation: Example

Predicates S (same generation), p (parent), h (human)

$$S(x, x) \leftarrow h(x)$$

$$S(x, y) \leftarrow p(x, w) \wedge S(w, y) \wedge p(y, w)$$

with query $S(1, x)$.

~> Query rule: $\text{Query}(x) \leftarrow S(1, x)$

Transformed rules:

$$\text{Query}^f(x) \leftarrow S^{bf}(1, x)$$

QSQR Transformation: Example

Predicates S (same generation), p (parent), h (human)

$$S(x, x) \leftarrow h(x)$$

$$S(x, y) \leftarrow p(x, w) \wedge S(w, y) \wedge p(y, w)$$

with query $S(1, x)$.

~> Query rule: $\text{Query}(x) \leftarrow S(1, x)$

Transformed rules:

$$\text{Query}^f(x) \leftarrow S^{bf}(1, x)$$

$$S^{bf}(x, x) \leftarrow h(x)$$

QSQR Transformation: Example

Predicates S (same generation), p (parent), h (human)

$$S(x, x) \leftarrow h(x)$$

$$S(x, y) \leftarrow p(x, w) \wedge S(w, y) \wedge p(y, y)$$

with query $S(1, x)$.

~> Query rule: $\text{Query}(x) \leftarrow S(1, x)$

Transformed rules:

$$\text{Query}^f(x) \leftarrow S^{bf}(1, x)$$

$$S^{bf}(x, x) \leftarrow h(x)$$

$$S^{bf}(x, y) \leftarrow p(x, w) \wedge S^{fb}(w, y) \wedge p(y, y)$$

QSQR Transformation: Example

Predicates S (same generation), p (parent), h (human)

$$S(x, x) \leftarrow h(x)$$

$$S(x, y) \leftarrow p(x, w) \wedge S(w, y) \wedge p(y, y)$$

with query $S(1, x)$.

~> Query rule: $\text{Query}(x) \leftarrow S(1, x)$

Transformed rules:

$$\text{Query}^f(x) \leftarrow S^{bf}(1, x)$$

$$S^{bf}(x, x) \leftarrow h(x)$$

$$S^{bf}(x, y) \leftarrow p(x, w) \wedge S^{fb}(w, y) \wedge p(y, y)$$

$$S^{fb}(x, x) \leftarrow h(x)$$

QSQR Transformation: Example

Predicates S (same generation), p (parent), h (human)

$$S(x, x) \leftarrow h(x)$$

$$S(x, y) \leftarrow p(x, w) \wedge S(w, w) \wedge p(y, w)$$

with query $S(1, x)$.

~> Query rule: $\text{Query}(x) \leftarrow S(1, x)$

Transformed rules:

$$\text{Query}^f(x) \leftarrow S^{bf}(1, x)$$

$$S^{bf}(x, x) \leftarrow h(x)$$

$$S^{bf}(x, y) \leftarrow p(x, w) \wedge S^{fb}(w, w) \wedge p(y, w)$$

$$S^{fb}(x, x) \leftarrow h(x)$$

$$S^{fb}(x, y) \leftarrow p(x, w) \wedge S^{fb}(w, w) \wedge p(y, w)$$

Summary and Outlook

Datalog queries can be evaluated bottom-up or top-down

Simplest practical bottom-up technique: semi-naive evaluation

Top-down: Query-Subquery (QSQ) approach (goal-directed)

Next question:

- Can bottom-up evaluations be goal directed?
- What about practical implementations?
- Graph databases