

Algorithmic Aspects of Semiring Provenance for Stratified Datalog

Matthias Naaf



Logic and Algebra for Query Evaluation, Berkeley 2023

Algorithmic Aspects

Semiring Provenance for Stratified Datalog

Computing Greatest Fixed Points
(in absorptive semirings)

Circuit Representations

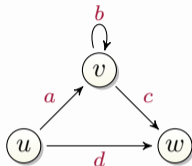
Why Greatest Fixed Points?

Semiring Semantics for Datalog

Datalog

$$Txy :- Exy$$

$$Txy :- Exz, Tzy$$



Equation System

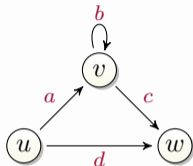
$$T_{uv} = a \vee (a \wedge T_{vv}) \vee (d \wedge T_{wv})$$

$$T_{uw} = d \vee (d \wedge T_{ww}) \vee (a \wedge T_{vw})$$

$$\vdots$$

Semiring Semantics for Datalog

Datalog

$$Txy :- Exy$$
$$Txy :- Exz, Tzy$$


Equation System

$$T_{uv} = a + (a \cdot T_{vv}) + (d \cdot T_{wv})$$

$$T_{uw} = d + (d \cdot T_{ww}) + (a \cdot T_{vw})$$

$$\vdots$$

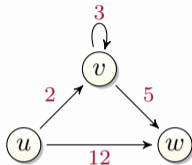
Semantics: Least solution

► Power series: $T_{uw}^* = d + ac + abc + ab^2c + ab^3c + \dots$

► PosBool: $T_{uw}^* = d \vee (a \wedge c)$

Semiring Semantics for Datalog

Datalog

$$Txy :- Exy$$
$$Txy :- Exz, Tzy$$


Equation System

$$T_{uv} = a + (a \cdot T_{vv}) + (d \cdot T_{wv})$$

$$T_{uw} = d + (d \cdot T_{ww}) + (a \cdot T_{vw})$$

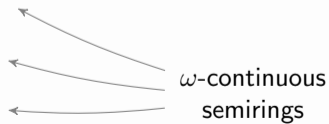
$$\vdots$$

Semantics: Least solution

▶ Power series: $T_{uw}^* = d + ac + abc + ab^2c + ab^3c + \dots$

▶ PosBool: $T_{uw}^* = d \vee (a \wedge c)$

▶ Tropical: $T_{uw}^* = \min(12, 2 + 5) = 7$



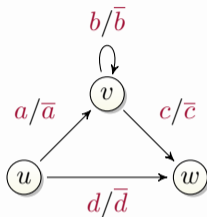
Semiring Semantics for Stratified Datalog

Stratified Datalog

$$Txy :- Exy$$

$$Txy :- Exz, Tzy$$

$$Nxy :- \neg Txy$$



Negation: can be defined in some semirings

► PosBool: $T_{uw}^* = d \vee (a \wedge c)$

$$N_{uw}^* = \overline{T_{uw}^*} = \bar{d} \wedge (\bar{a} \vee \bar{c})$$

but not clear how do to it in general

► Polynomials: $\overline{a^2} = ?$

► Tropical: $\bar{7} = ?$

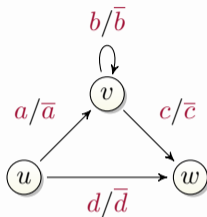
Semiring Semantics for Stratified Datalog

Stratified Datalog

$$Txy := Exy$$

$$Txy := Exz, Tzy$$

$$Nxy := \neg Txy$$



Equation System

$$T_{uv} = a + (a \cdot T_{vv} + d \cdot T_{wv})$$

$$T_{uw} = d + (d \cdot T_{ww} + a \cdot T_{vw})$$

\implies Least solution

Dualized System

$$N_{uv} = \bar{a} \cdot (\bar{a} + N_{vv}) \cdot (\bar{d} + N_{wv})$$

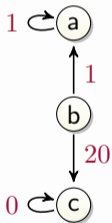
$$N_{uw} = \bar{d} \cdot (\bar{d} + N_{ww}) \cdot (\bar{a} + N_{vw})$$

\implies Greatest solution

Motivation II: Fixed-point Logic

$$[\mathbf{gfp} \ Rx. \exists y (Exy \wedge Ry)](v)$$

“there is an infinite path from v ”



$$R_a = 1 + R_a$$

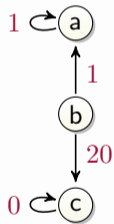
$$R_b = \min(1 + R_a, 20 + R_c)$$

$$R_c = 0 + R_c$$

Motivation II: Fixed-point Logic

$$[\mathbf{gfp} \ Rx. \exists y (Exy \wedge Ry)](v)$$

cost of
“~~there is~~ an infinite path from v ”



$$R_a = 1 + R_a$$

$$R_b = \min(1 + R_a, 20 + R_c)$$

$$R_c = 0 + R_c$$

$$R_a^* = \infty$$

$$R_b^* = 20$$

$$R_c^* = 0$$

Greatest Solution

Computing Greatest Fixed Points

Naive Approach

$$\begin{aligned}R_a &= 1 + R_a \\R_b &= \min(1 + R_a, 20 + R_c) \\R_c &= 0 + R_c\end{aligned}$$

Naive Approach

Goal: Compute greatest fixed point of a polynomial operator

$$\mathbf{F} : \begin{pmatrix} R_a \\ R_b \\ R_c \end{pmatrix} \mapsto \begin{pmatrix} 1 + R_a \\ \min(1 + R_a, 20 + R_c) \\ 0 + R_c \end{pmatrix}$$



Naive Approach

Goal: Compute greatest fixed point of a polynomial operator

$$\mathbf{F} : \begin{pmatrix} R_a \\ R_b \\ R_c \end{pmatrix} \mapsto \begin{pmatrix} 1 + R_a \\ \min(1 + R_a, 20 + R_c) \\ 0 + R_c \end{pmatrix}$$

Iteration:

$$\begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} \mapsto \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix} \mapsto \begin{pmatrix} 2 \\ 2 \\ 0 \end{pmatrix} \mapsto \dots \mapsto \begin{pmatrix} 20 \\ 20 \\ 0 \end{pmatrix} \mapsto \begin{pmatrix} 21 \\ 20 \\ 0 \end{pmatrix} \mapsto \begin{pmatrix} 22 \\ 20 \\ 0 \end{pmatrix} \mapsto \dots \mapsto \begin{pmatrix} \infty \\ 20 \\ 0 \end{pmatrix}$$

Main Result

Let $(K, +, \cdot, 0, 1)$ be an absorptive, fully-continuous semiring.
For a polynomial operator $\mathbf{F}: K^n \rightarrow K^n$,

$$\text{lfp}(\mathbf{F}) = \mathbf{F}^n(\mathbf{0}), \quad \text{gfp}(\mathbf{F}) = \mathbf{F}^n(\mathbf{F}^n(\mathbf{1})^\infty).$$

We only need a **polynomial number** of semiring operations:

$$\underbrace{\begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} \mapsto \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix} \mapsto \begin{pmatrix} 2 \\ 2 \\ 0 \end{pmatrix} \mapsto \begin{pmatrix} 3 \\ 3 \\ 0 \end{pmatrix}}_{\leq n} \xrightarrow{\infty} \underbrace{\begin{pmatrix} \infty \\ \infty \\ 0 \end{pmatrix} \mapsto \begin{pmatrix} \infty \\ 20 \\ 0 \end{pmatrix}}_{\leq n} \curvearrowright$$

① Fully continuous

- ▶ Natural order: $a \leq a + b$
- ▶ Each chain has supremum $\bigsqcup C$ and infimum $\bigsqcap C$, these commute with $+/\cdot$

② Absorption

- ▶ $a + a \cdot b = a \iff 1 \text{ is greatest element} \iff a \cdot b \leq a$

1 Fully continuous

- ▶ Natural order: $a \leq a + b$
- ▶ Each chain has supremum $\bigsqcup C$ and infimum $\bigsqcap C$, these commute with $+/\cdot$

2 Absorption

- ▶ $a + a \cdot b = a \iff 1 \text{ is greatest element} \iff a \cdot b \leq a$

Infinitary Power

For $a \in K$ we define: $a^\infty := \bigsqcap_{n < \omega} a^n$



Remember:

Decreasing multiplication

Main Result

Let $(K, +, \cdot, 0, 1)$ be an absorptive, fully-continuous semiring.
For a polynomial operator $\mathbf{F}: K^n \rightarrow K^n$,

$$\text{lfp}(\mathbf{F}) = \mathbf{F}^n(\mathbf{0}), \quad \text{gfp}(\mathbf{F}) = \mathbf{F}^n(\mathbf{F}^n(\mathbf{1})^\infty).$$

Proof sketch:



derivation trees

+



absorption

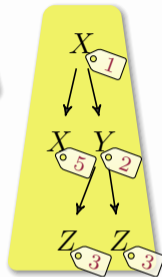
Derivation Trees

$$\begin{pmatrix} X \\ Y \\ Z \end{pmatrix} \mapsto \begin{pmatrix} \min(5, 1+X+Y) \\ 2+Z+Z \\ \min(3, 1+Z) \end{pmatrix}$$

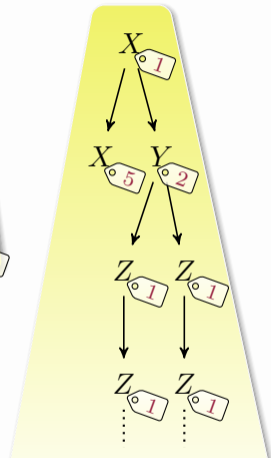
inspired by Newton's method
(Esparza, Kiefer, Luttenberger, JACM'10)



cost: 5



cost: 14



cost: $8 + 2 + 2 + \dots = \infty$

Derivation Trees

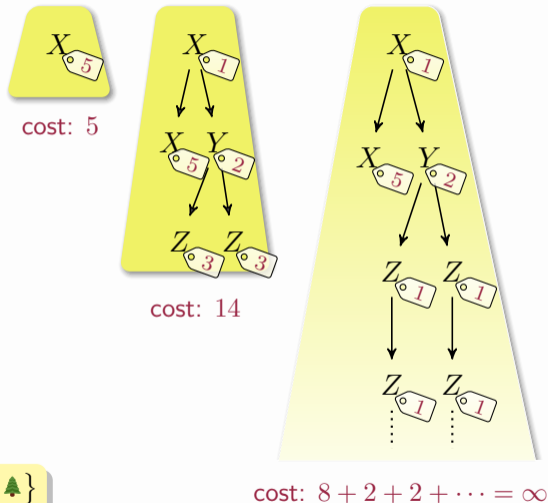
$$\begin{pmatrix} X \\ Y \\ Z \end{pmatrix} \mapsto \begin{pmatrix} \min(5, 1+X+Y) \\ 2+Z+Z \\ \min(3, 1+Z) \end{pmatrix}$$

inspired by Newton's method
(Esparza, Kiefer, Luttenberger, JACM'10)



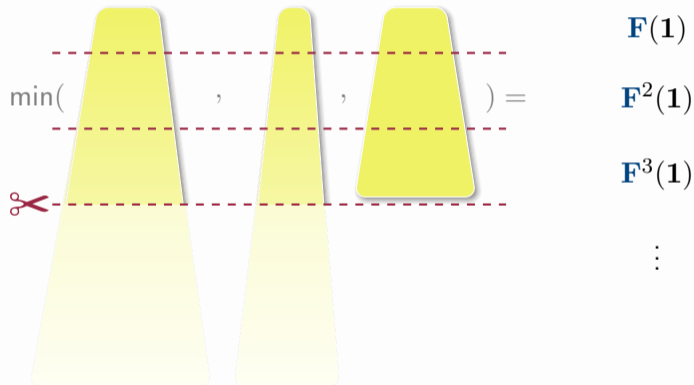
$\text{lfp} = \min \{ \text{cost}(\text{tree}) \mid \text{finite tree} \}$

$\text{gfp} = \min \{ \text{cost}(\text{tree}) \mid \text{finite tree, infinite tree} \}$



Derivation Trees vs. Iteration

Observation: Prefixes of  correspond to iteration steps.

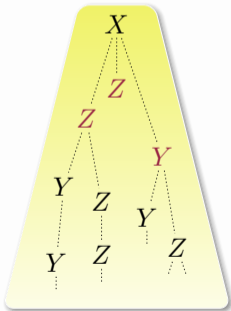


$$\bigsqcap_{n < \omega} : \min \left\{ \text{cost}(\text{tree}) \mid \text{finite/infinite tree} \right\} = \text{gfp}(\mathbf{F}) \quad \blacksquare$$

Absorption on Derivation Trees



If each coefficient 2 occurs more often in \bullet than in \bullet , then $\text{cost}(\bullet)$ is **absorbed by** $\text{cost}(\bullet)$.



complicated tree \bullet

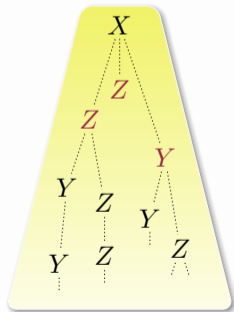


nice tree \bullet

Absorption on Derivation Trees

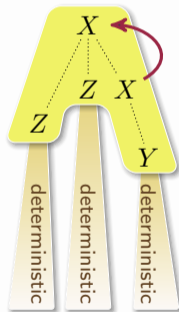


If each coefficient $\langle 2 \rangle$ occurs more often in \bullet than in \bullet , then $\text{cost}(\bullet)$ is **absorbed by** $\text{cost}(\bullet)$.



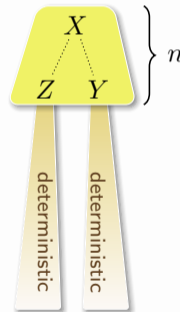
complicated tree \bullet

\geq
cost



ultimately periodic

\geq
cost

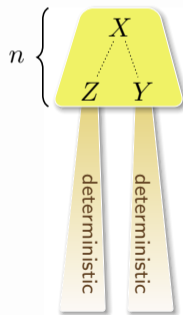


nice tree \bullet

Computing Nice Trees

Main Result

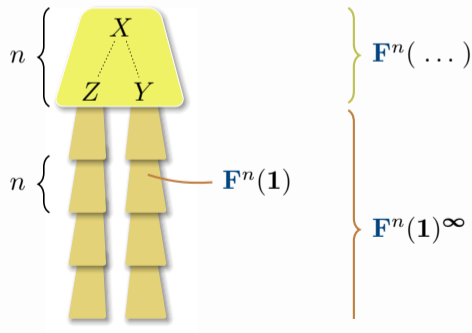
$$\text{gfp}(\mathbf{F}) = \min \left\{ \text{cost}(\text{🌲}) \mid \text{nice } \text{🌲} \right\} = \dots$$



Computing Nice Trees

Main Result

$$\text{gfp}(\mathbf{F}) = \min \left\{ \text{cost}(\text{🌲}) \mid \text{nice } \text{🌲} \right\} = \mathbf{F}^n(\mathbf{F}^n(\mathbf{1})^\infty)$$



Back to Datalog: Circuits

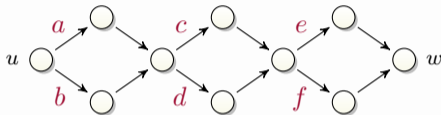
Circuits for Datalog Provenance

Problem: Provenance in polynomial semirings can become large

Datalog

$$Txy :- Exy$$

$$Txy :- Exz, Tzy$$



$$\text{PosBool: } T_{uw}^* = ace + acf + ade + adf + bce + bcf + bde + bdf$$

Solution: Represent provenance computation by a small circuit

Circuits for Datalog Provenance

Recall

$$\text{lfp}(\mathbf{F}) = \mathbf{F}^n(\mathbf{0})$$

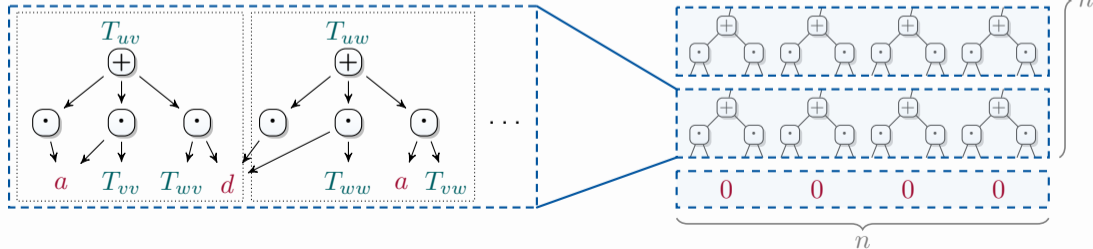
Equation System

$$T_{uv} = a + (a \cdot T_{vv}) + (d \cdot T_{wv})$$

$$T_{uw} = d + (d \cdot T_{ww}) + (a \cdot T_{vw})$$

⋮

⋮



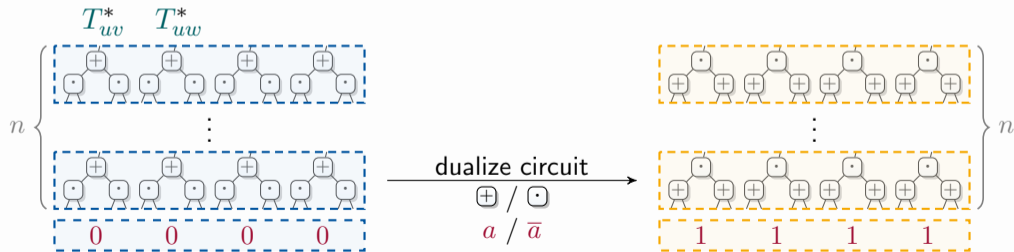
Circuits for Stratified Datalog

Strat. Datalog

$$Txy :- Exy$$

$$Txy :- Exz, Tzy$$

$$Nxy :- \neg Txy$$



Circuits for Stratified Datalog

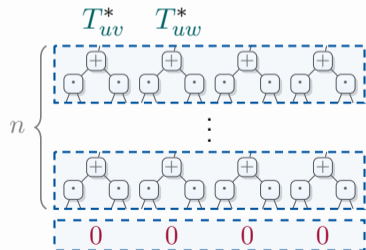
$$\text{gfp}(\mathbf{F}) = \mathbf{F}^n(\mathbf{F}^n(\mathbf{1})^\infty)$$

Strat. Datalog

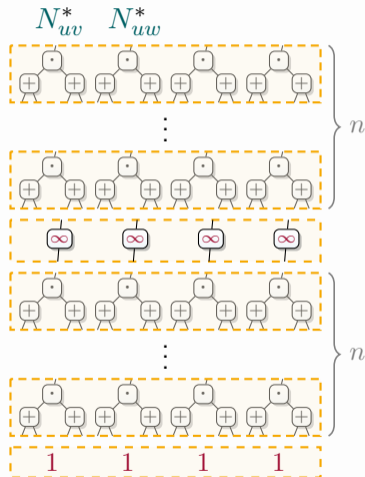
$$Txy := Exy$$

$$Txy := Exz, Tzy$$

$$Nxy := \neg Txy$$



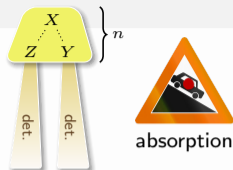
dualize circuit
 \oplus / \ominus
 a / \bar{a}



Summary

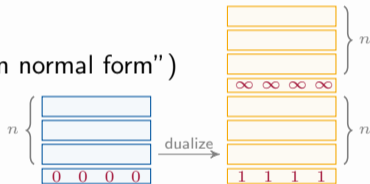
Computing greatest fixed points

- ▶ In absorptive semirings: $\text{gfp}(\mathbf{F}) = \mathbf{F}^n(\mathbf{F}^n(\mathbf{1})^\infty)$



Semiring provenance for stratified Datalog

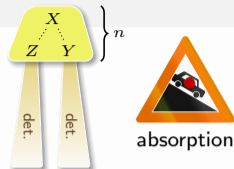
- ▶ Negation: **greatest** solution to **dual equation system** (“negation normal form”)
- ▶ Circuit representations for Datalog can be generalized



Summary

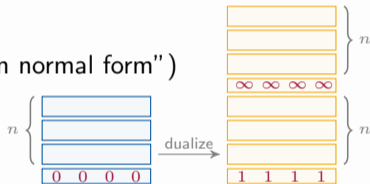
Computing greatest fixed points

- ▶ In absorptive semirings: $\text{gfp}(\mathbf{F}) = \mathbf{F}^n (\mathbf{F}^n (\mathbf{1})^\infty)$



Semiring provenance for stratified Datalog

- ▶ Negation: **greatest** solution to **dual equation system** (“negation normal form”)
- ▶ Circuit representations for Datalog can be generalized



Questions

1 Applications

- ▶ LFP: strategies in infinite games
- ▶ Stratified Datalog: ?

2 Alternating fixed points

- ▶ Is the main result applicable?
- ▶ Quasipolynomial time?