

Tuple-Generating Dependencies Capture Complex Values

Maximilian Marx, Markus Krötzsch

Knowledge-based systems group
TU Dresden

ICDT 2022

Full paper: <https://iccl.inf.tu-dresden.de/web/DataLogCV/en>

Datalog^{CV}: Datalog with Complex Values

Overview:

- ▶ Datalog^{CV} extends Datalog with sorts for **tuples** and **sets**
- ▶ these sorts can be nested: e.g., $\{\langle\{0, 1\}, 1\rangle\}$
- ▶ sets support unions $S \cup T$ and intersections $S \cap T$
- ▶ like Datalog, entailment $\mathbb{P}, \mathbb{D} \models \alpha$ can be decided via **least model** of \mathbb{P} and \mathbb{D}
- ▶ databases are **flat**, i.e., don't contain sets or tuples

Datalog^{CV}: Datalog with Complex Values

Overview:

- ▶ Datalog^{CV} extends Datalog with sorts for **tuples** and **sets**
- ▶ these sorts can be nested: e.g., $\{\langle\{0, 1\}, 1\rangle\}$
- ▶ sets support unions $S \cup T$ and intersections $S \cap T$
- ▶ like Datalog, entailment $\mathbb{P}, \mathbb{D} \models \alpha$ can be decided via **least model** of \mathbb{P} and \mathbb{D}
- ▶ databases are **flat**, i.e., don't contain sets or tuples

Features:

- ▶ the usual set operations $x \in S, x \notin S, S \subseteq T, S \subset T$, and $\mathcal{P}(S)$ are definable
- ▶ inequality is axiomatisable for any sort

Datalog^{CV}: Datalog with Complex Values

Overview:

- ▶ Datalog^{CV} extends Datalog with sorts for **tuples** and **sets**
- ▶ these sorts can be nested: e.g., $\{\langle\{0, 1\}, 1\rangle\}$
- ▶ sets support unions $S \cup T$ and intersections $S \cap T$
- ▶ like Datalog, entailment $\mathbb{P}, \mathbb{D} \models \alpha$ can be decided via **least model** of \mathbb{P} and \mathbb{D}
- ▶ databases are **flat**, i.e., don't contain sets or tuples

Features:

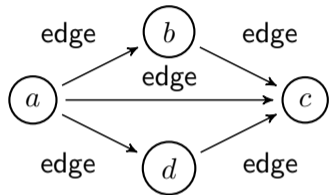
- ▶ the usual set operations $x \in S, x \notin S, S \subseteq T, S \subset T$, and $\mathcal{P}(S)$ are definable
- ▶ inequality is axiomatisable for any sort

Translation:

- ▶ Datalog^{CV} programs admit a translation into sets of tuple-generating dependencies
- ▶ translated programs are all-instances all-strategies terminating under the standard chase
- ▶ use nulls to reify sets and tuples, represent sets as unions of singleton sets

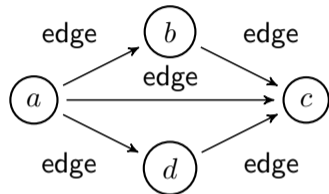
Datalog^{CV} by Example

- ▶ Consider a database encoding a graph as $\text{edge}(x, y)$:



Datalog^{CV} by Example

- ▶ Consider a database encoding a graph as $\text{edge}(x, y)$:

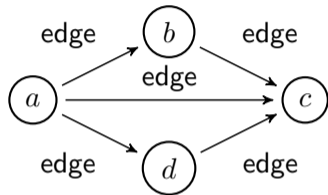


- ▶ Query for all paths from x to y :

$$\begin{aligned} \text{edge}(x, y) &\rightarrow \text{path}(x, y, \{\langle x, y \rangle\}) \\ \text{path}(x, y, P) \wedge \text{edge}(y, z) &\rightarrow \text{path}(x, z, P \cup \{\langle y, z \rangle\}) \end{aligned}$$

Datalog^{CV} by Example

- ▶ Consider a database encoding a graph as $\text{edge}(x, y)$:



- ▶ Query for all paths from x to y :

$$\begin{aligned} \text{edge}(x, y) &\rightarrow \text{path}(x, y, \{\langle x, y \rangle\}) \\ \text{path}(x, y, P) \wedge \text{edge}(y, z) &\rightarrow \text{path}(x, z, P \cup \{\langle y, z \rangle\}) \end{aligned}$$

- ▶ Paths:

$$\begin{array}{lll} \text{path}(a, c, \{\langle a, c \rangle\}) & \text{path}(a, b, \{\langle a, b \rangle\}) & \text{path}(a, d, \{\langle a, d \rangle\}) \\ \text{path}(a, c, \{\langle a, b \rangle, \langle b, c \rangle\}) & \text{path}(a, c, \{\langle a, d \rangle, \langle d, c \rangle\}) & \end{array}$$

Schema Heights:

set-height maximal nesting depth of set sorts

tuple-height maximal nesting depth of tuple sorts

Theorem

Fact entailment for Datalog^{CV} with set-height k is

- ▶ *k EXPTIME-complete for data complexity,*
- ▶ *$(k + 2)$ EXPTIME-complete for combined complexity, and*

Schema Heights:

set-height maximal nesting depth of set sorts

tuple-height maximal nesting depth of tuple sorts

Theorem

Fact entailment for Datalog^{CV} with set-height k is

- ▶ *k EXPTIME-complete for data complexity,*
- ▶ *$(k + 2)$ EXPTIME-complete for combined complexity, and*
- ▶ *$(k + 1)$ EXPTIME-complete for combined complexity if the tuple-height is bounded.*

Schema Heights:

set-height maximal nesting depth of set sorts

tuple-height maximal nesting depth of tuple sorts

Theorem

Fact entailment for Datalog^{CV} with set-height k is

- ▶ k EXPTIME-complete for data complexity,
- ▶ $(k + 2)$ EXPTIME-complete for combined complexity, and
- ▶ $(k + 1)$ EXPTIME-complete for combined complexity if the tuple-height is bounded.

- ▶ Lower bound by extending linear orders to powersets
- ▶ Upper bound from translation into TGDs

A Tractable Fragment

- ▶ Is there a non-trivial fragment with PTIME data complexity?
- ▶ Observation: if all sets are bounded, we can replace sets by tuples
- ▶ \rightsquigarrow bounded cardinality Datalog^{CV} is essentially high-arity Datalog

A Tractable Fragment

- ▶ Is there a non-trivial fragment with PTIME data complexity?
- ▶ Observation: if all sets are bounded, we can replace sets by tuples
- ▶ \rightsquigarrow bounded cardinality Datalog^{CV} is essentially high-arity Datalog

Definition

- ▶ ground fact α has k -bounded cardinality if all set terms occurring in α are of the form $\{t_1, \dots, t_n\}$ with $n \leq k$
- ▶ program \mathbb{P} has k -bounded cardinality if, for all databases \mathbb{D} , all ground facts α with $\mathbb{P}, \mathbb{D} \models \alpha$ have k -bounded cardinality
- ▶ \mathbb{P} has bounded cardinality if it has k -bounded cardinality for some k

A Tractable Fragment

- ▶ Is there a non-trivial fragment with PTIME data complexity?
- ▶ Observation: if all sets are bounded, we can replace sets by tuples
- ▶ \rightsquigarrow bounded cardinality Datalog^{CV} is essentially high-arity Datalog

Definition

- ▶ ground fact α has k -bounded cardinality if all set terms occurring in α are of the form $\{t_1, \dots, t_n\}$ with $n \leq k$
- ▶ program \mathbb{P} has k -bounded cardinality if, for all databases \mathbb{D} , all ground facts α with $\mathbb{P}, \mathbb{D} \models \alpha$ have k -bounded cardinality
- ▶ \mathbb{P} has bounded cardinality if it has k -bounded cardinality for some k

Example:

$$\begin{array}{l} e(x) \rightarrow s(\{x\}) \quad \text{has 2-bounded cardinality.} \\ s(X) \wedge s(Y) \rightarrow p(X \cup Y) \end{array}$$

A Tractable Fragment

- ▶ Is there a non-trivial fragment with PTIME data complexity?
- ▶ Observation: if all sets are bounded, we can replace sets by tuples
- ▶ \rightsquigarrow bounded cardinality Datalog^{CV} is essentially high-arity Datalog

Definition

- ▶ ground fact α has k -bounded cardinality if all set terms occurring in α are of the form $\{t_1, \dots, t_n\}$ with $n \leq k$
- ▶ program \mathbb{P} has k -bounded cardinality if, for all databases \mathbb{D} , all ground facts α with $\mathbb{P}, \mathbb{D} \models \alpha$ have k -bounded cardinality
- ▶ \mathbb{P} has bounded cardinality if it has k -bounded cardinality for some k

Example:

$$\begin{array}{ll} e(x) \rightarrow s(\{x\}) & \text{does not have bounded cardinality.} \\ s(X) \wedge s(Y) \rightarrow s(X \cup Y) & \end{array}$$

Bounded Cardinality: Bad News & Good News

Bad News:

Theorem

It is undecidable whether a Datalog^{CV} program has bounded cardinality.

Bounded Cardinality: Bad News & Good News

Bad News:

Theorem

It is undecidable whether a Datalog^{CV} program has bounded cardinality.

Good News:

Theorem

For fixed k , deciding if a Datalog^{CV} program \mathbb{P} has k -bounded cardinality is 2EXPTIME-complete (EXPTIME-complete if set-height and tuple-height are bounded).

Theorem

Bounded-cardinality Datalog^{CV} is PTIME-complete for data complexity and 2EXPTIME-complete for combined complexity.

A Syntactic Criterion: Weak Set-Acyclicity

- ▶ 'large' sets are constructed recursively as unions of smaller sets
- ▶ idea: track propagation of sets along positions of a program
- ▶ if there are no cycles involving unions, sets must be bounded
- ▶ like Weak Acyclicity, this can be checked using a propagation graph

A Syntactic Criterion: Weak Set-Acyclicity

- ▶ 'large' sets are constructed recursively as unions of smaller sets
- ▶ idea: track propagation of sets along positions of a program
- ▶ if there are no cycles involving unions, sets must be bounded
- ▶ like Weak Acyclicity, this can be checked using a propagation graph

Example

$$e(x) \rightarrow s(\{x\})$$

$$s(X) \wedge s(Y) \rightarrow p(X \cup Y)$$

has WSA graph

$$s \cdot \epsilon \Rightarrow p \cdot \epsilon$$

$$e \cdot \epsilon \rightarrow s \cdot 1 \rightarrow p \cdot 1$$

no cycle along 'special' edges
 \rightsquigarrow WSA

A Syntactic Criterion: Weak Set-Acyclicity

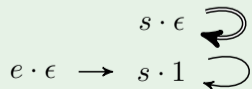
- ▶ 'large' sets are constructed recursively as unions of smaller sets
- ▶ idea: track propagation of sets along positions of a program
- ▶ if there are no cycles involving unions, sets must be bounded
- ▶ like Weak Acyclicity, this can be checked using a propagation graph

Example

$$e(x) \rightarrow s(\{x\})$$

$$s(X) \wedge s(Y) \rightarrow s(X \cup Y)$$

has WSA graph



cycle along 'special' edge $s \cdot \epsilon \Rightarrow s \cdot \epsilon$
 \rightsquigarrow not WSA

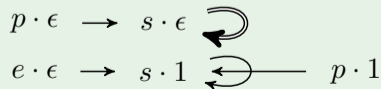
A Syntactic Criterion: Weak Set-Acyclicity

- ▶ 'large' sets are constructed recursively as unions of smaller sets
- ▶ idea: track propagation of sets along positions of a program
- ▶ if there are no cycles involving unions, sets must be bounded
- ▶ like Weak Acyclicity, this can be checked using a propagation graph

Example

$$e(x) \rightarrow s(\{x\}) \quad s(X) \wedge s(Y) \wedge p(P) \rightarrow s(P \cap (X \cup Y)) \quad s(X) \wedge s(Y) \rightarrow p(X \cup Y)$$

has WSA graph



cycle along 'special' edge $s \cdot \epsilon \Rightarrow s \cdot \epsilon$

\rightsquigarrow not WSA (but 2-bounded cardinality)

Cardinality Constraints

- ▶ For every rule, bound the maximal cardinality of sets derived (w.r.t. sets in the match) \rightsquigarrow system of inequalities
- ▶ If the sum of the bounds w.r.t. the inequalities has a minimum, the program has bounded cardinality

Cardinality Constraints

- ▶ For every rule, bound the maximal cardinality of sets derived (w.r.t. sets in the match) \rightsquigarrow system of inequalities
- ▶ If the sum of the bounds w.r.t. the inequalities has a minimum, the program has bounded cardinality

Example

$$e(x) \rightarrow s(\{x\}) \quad s(X) \wedge s(Y) \wedge p(S) \rightarrow s(S \cap (X \cup Y)) \quad s(X) \wedge s(Y) \rightarrow p(X \cup Y)$$

has cardinality constraints

$$x_{s \cdot \epsilon} \geq 0 \quad x_{p \cdot \epsilon} \geq 0 \quad x_{s \cdot \epsilon} \geq 1 \quad x_{p \cdot \epsilon} \geq 2 \quad x_{s \cdot \epsilon} \geq \min(x_{p \cdot \epsilon}, x_{s \cdot \epsilon} + x_{s \cdot \epsilon})$$

with optimal solution $x_{s \cdot \epsilon} = 2 = x_{p \cdot \epsilon}$ and optimal value $k = 4 \rightsquigarrow$ bounded cardinality

Conclusions:

- ▶ Datalog^{CV} is a highly expressive query language
- ▶ TGD translation offers a path to reasoning with existing software
- ▶ Bounded cardinality Datalog^{CV} is a tractable fragment
- ▶ Weak Set-Acyclicity and Cardinality Constraints are sufficient conditions for bounded cardinality

Outlook:

- ▶ Can Datalog^{CV} express all decidable monotonic queries?
- ▶ Are there other tractable fragments?
- ▶ What happens if we add stratified negation?

Powerset:

$$\rightarrow \text{PSU}_\tau(x, \emptyset, \emptyset)$$

$$\text{PSU}_\tau(x, P, Q) \rightarrow \text{PSU}_\tau(x, P \cup \{S\}, Q \cup \{S \cup \{x\}\})$$

$$\rightarrow \text{PS}_\sigma(\emptyset, \{\emptyset\})$$

$$\text{PS}_\sigma(S, P) \wedge \text{PSU}_\tau(x, P, Q) \rightarrow \text{PS}_\sigma(S \cup \{x\}, P \cup Q)$$

Subsets, Containment, Inequality:

$$\rightarrow S \subseteq_\sigma S \cup T$$

$$S \cup \{x\} \subseteq_\sigma S \rightarrow x \in_\sigma S$$

$$S \cap \{x\} \subseteq_\sigma \emptyset \rightarrow x \notin_\sigma S$$

$$\{x\} \cap \{y\} \subseteq_{\{\tau\}} \emptyset \rightarrow x \neq_\tau y$$

$$S \subseteq_\sigma T \wedge x \in_\sigma T \wedge x \notin_\sigma S \rightarrow S \subset_\sigma T$$

$$\begin{aligned} \text{first}_\tau(x) \wedge \text{last}_\tau(z) &\rightarrow \text{first}_\sigma(\emptyset_\sigma) \wedge \\ &\quad \text{mkStep}_\sigma(\emptyset_\sigma, x, x) \wedge \\ &\quad \text{last}_\sigma(\{z\}) \end{aligned}$$

$$\text{mkStep}_\sigma(S, t, t) \rightarrow \text{next}_\sigma(S, S \cup \{t\})$$

$$\text{mkStep}_\sigma(S, t, x) \wedge \text{last}_\tau(x) \rightarrow \text{step}_\sigma(S, t, x, S \cup \{x\})$$

$$\begin{aligned} \text{mkStep}_\sigma(S, t, x) \wedge \text{next}_\tau(x, y) &\rightarrow \text{mkStep}_\sigma(S \cup \{x\}, y, y) \wedge \\ &\quad \text{mkStep}_\sigma(S, t, y) \end{aligned}$$

$$\begin{aligned} \text{mkStep}_\sigma(S, t, x) \wedge \text{next}_\tau(x, y) \wedge \\ \text{step}_\sigma(S \cup \{x\}, y, y, X) \wedge \text{step}_\sigma(S, t, y, Z) &\rightarrow \text{next}_\sigma(X, S \cup \{y\}) \wedge \\ &\quad \text{step}_\sigma(S, t, x, Z) \end{aligned}$$

Tuple sorts:

$$\bigwedge_{i=1}^{\ell} \text{sort}_{\tau_i}(x_i) \rightarrow \exists z. \text{tuple}_{\pi}(z, x_1, \dots, x_{\ell}) \wedge \text{sort}_{\pi}(z)$$

Set sorts:

$$\rightarrow \exists V. \text{empty}_{\sigma}(V) \wedge \text{sort}_{\sigma}(V) \wedge \text{done}_{\sigma}(V)$$

$$\text{done}_{\sigma}(V) \wedge \text{sort}_{\tau}(x) \rightarrow \exists W. \text{SU}_{\sigma}(x, V, W) \wedge \text{sort}_{\sigma}(W)$$

$$\wedge \text{todo}_{\sigma}(W, W)$$

$$\text{todo}_{\sigma}(V, W) \wedge \text{SU}_{\sigma}(x, U, V) \rightarrow \text{SU}_{\sigma}(x, W, W) \wedge \text{todo}_{\sigma}(U, W)$$

$$\text{todo}_{\sigma}(V, W) \wedge \text{empty}_{\sigma}(V) \rightarrow \text{done}_{\sigma}(W)$$

Subsets, Unions, Not-in:

$$\begin{aligned} & \text{sort}_\sigma(V) \wedge \text{sort}_\sigma(W) \rightarrow \text{ckSub}_\sigma(V, V, W) \\ \text{ckSub}_\sigma(U, V, W) \wedge \text{SU}_\sigma(x, U', U) \wedge \text{SU}_\sigma(x, V, V) & \rightarrow \text{ckSub}_\sigma(U', V, W) \\ & \text{ckSub}_\sigma(U, V, W) \wedge \text{empty}_\sigma(U) \rightarrow \text{subset}_\sigma(V, W) \\ & \text{subset}_\sigma(V, W) \wedge \text{subset}_\sigma(W, V) \rightarrow \text{eq}_\sigma(V, W) \\ & \text{empty}_\sigma(V) \wedge \text{sort}_\sigma(W) \rightarrow \text{U}_\sigma(V, W, W) \\ \text{SU}_\sigma(x, W, W') \wedge \text{U}_\sigma(V, W', U) \wedge \text{SU}_\sigma(x, V, V') & \rightarrow \text{U}_\sigma(V', W, U) \\ & \text{sort}_\sigma(V) \wedge \text{sort}_\tau(x) \rightarrow \text{ckNIn}_\sigma(V, x, V) \\ \text{ckNIn}_\sigma(U, x, V) \wedge \text{SU}_\sigma(y, U', U) \wedge \text{NEq}_\tau(x, y) & \rightarrow \text{ckNIn}_\sigma(U', x, V) \\ & \text{ckNIn}_\sigma(U, x, V) \wedge \text{empty}_\sigma(U) \rightarrow \text{NIn}_\sigma(x, V) \end{aligned}$$

Tuples, Inequality, Intersections:

$$\text{tuple}_\pi(z, x_1, \dots, x_\ell) \wedge \text{tuple}_\pi(z', y_1, \dots, y_\ell) \wedge \text{NEq}_{\tau_i}(x_i, y_i) \\ \rightarrow \text{NEq}_\pi(z, z')$$

$$\text{SU}_\sigma(x, V, V) \wedge \text{NIn}_\sigma(x, W) \rightarrow \text{NEq}_\sigma(V, W) \wedge \text{NEq}_\sigma(W, V)$$

$$\text{empty}_\sigma(V) \wedge \text{sort}_\sigma(W) \rightarrow \text{I}_\sigma(V, W, V)$$

$$\text{I}_\sigma(U, V, W) \wedge \text{SU}_\sigma(x, U, U') \wedge \text{NIn}_\sigma(x, V) \rightarrow \text{I}_\sigma(U', V, W)$$

$$\text{I}_\sigma(U, V, W) \wedge \text{SU}_\sigma(x, U, U') \wedge \text{SU}_\sigma(x, V, V') \\ \wedge \text{SU}_\sigma(x, W, W') \rightarrow \text{I}_\sigma(U', V', W')$$

Linear Datalog^{CV} is still intractable

Linear Datalog^{CV}:

- ▶ rules $\varphi \rightarrow \psi$ where both φ and ψ are single-atom
- ▶ we need to allow non-flat databases

Theorem

Linear Datalog^{CV} with set-height $k > 0$ is $(k - 1)\text{EXPTIME}$ -hard (data complexity) and $(k + 1)\text{EXPTIME}$ -hard (combined complexity; $k\text{EXPTIME}$ -hard for bounded tuple-height).

Idea:

- ▶ encode each predicate p as a position of sets of $\text{ar}(p)$ -tuples in a single facts predicate
- ▶ a rule $\rho = \varphi \rightarrow \psi$ turns into

$$\begin{aligned} & \text{facts}(y_1 \cup \text{ts}_\varphi(p_1), y_2 \cup \text{ts}_\varphi(p_2), \dots) \rightarrow \\ & \text{facts}(y_1 \cup \text{ts}_\rho(p_1), y_2 \cup \text{ts}_\rho(p_2), \dots) \end{aligned}$$

with $\text{ts}_\Phi(p)$ the set of terms t for which $p(t)$ occurs in Φ