

Datalog as a Query Language for Data Exchange Systems

Marcelo Arenas

PUC Chile

Joint work with Pablo Barceló and Juan Reutter

Outline

- ▶ The data exchange scenario
 - ▶ Query answering
- ▶ Queries with negation in data exchange
- ▶ $\text{DATALOG}^{\mathbf{C}(\neq)}$ programs
- ▶ Beyond union of conjunctive queries
 - ▶ Expressive power of $\text{DATALOG}^{\mathbf{C}(\neq)}$
 - ▶ New tractable classes of queries
- ▶ Concluding remarks

Outline

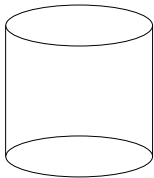
- ▶ The data exchange scenario
 - ▶ Query answering
- ▶ Queries with negation in data exchange
- ▶ $\text{DATALOG}^{\mathbf{C}(\neq)}$ programs
- ▶ Beyond union of conjunctive queries
 - ▶ Expressive power of $\text{DATALOG}^{\mathbf{C}(\neq)}$
 - ▶ New tractable classes of queries
- ▶ Concluding remarks

The data exchange scenario



The data exchange scenario

Source Schema **S** $\xrightarrow{\Sigma}$ Target Schema **T**



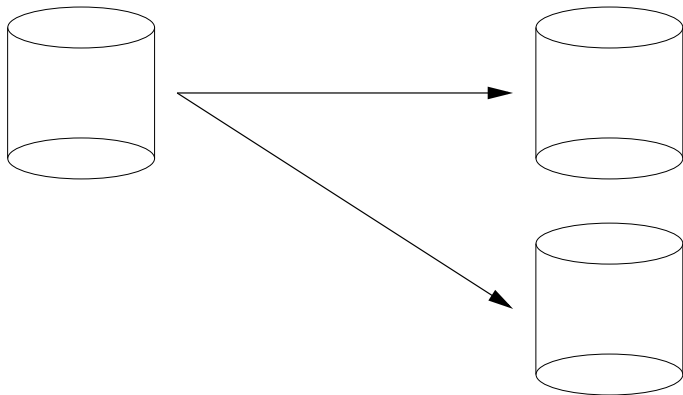
The data exchange scenario

Source Schema **S** $\xrightarrow{\Sigma}$ Target Schema **T**



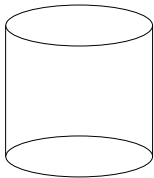
The data exchange scenario

Source Schema **S** $\xrightarrow{\Sigma}$ Target Schema **T**



The data exchange scenario

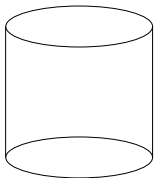
Source Schema **S** $\xrightarrow{\Sigma}$ Target Schema **T**



Query **Q**

The data exchange scenario

Source Schema **S** $\xrightarrow{\Sigma}$ Target Schema **T**

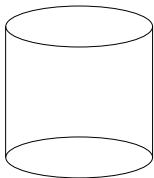


Query **Q**

What is the semantics of **Q**?

The data exchange scenario

Source Schema **S** $\xrightarrow{\Sigma}$ Target Schema **T**



Query **Q**

What is the semantics of **Q**?

Can **Q** be evaluated efficiently?

Data exchange settings

Data exchange setting: $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$

- ▶ **S**: source schema
- ▶ **T**: target schema
- ▶ Σ : set of **source-to-target dependencies**

Data exchange settings

Data exchange setting: $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$

- ▶ **S**: source schema
- ▶ **T**: target schema
- ▶ Σ : set of **source-to-target dependencies**

Source-to-target dependency:

$$\forall \bar{x} \forall \bar{y} \left(\varphi(\bar{x}, \bar{y}) \rightarrow \exists \bar{z} \psi(\bar{x}, \bar{z}) \right)$$

$\varphi(\bar{x}, \bar{y})$: conjunction of relational atomic formulas over **S**

$\psi(\bar{x}, \bar{z})$: conjunction of relational atomic formulas over **T**

Example: Data exchange setting

S: *Book*(*Title*, *AuthorName*, *Affiliation*)

T: *Writer*(*Name*, *BookTitle*, *Year*)

Σ :

$$Book(x_1, x_2, y_1) \rightarrow \exists z_1 Writer(x_2, x_1, z_1)$$

Data exchange problem

Given a source instance I , find a target instance J such that (I, J) satisfies Σ .

- ▶ (I, J) satisfies $\varphi(\bar{x}, \bar{y}) \rightarrow \exists \bar{z} \psi(\bar{x}, \bar{z})$ if for every (\bar{a}, \bar{b}) such that I satisfies $\varphi(\bar{a}, \bar{b})$, there is a tuple \bar{c} such that J satisfies $\psi(\bar{a}, \bar{c})$.

Data exchange problem

Given a source instance I , find a target instance J such that (I, J) satisfies Σ .

- ▶ (I, J) satisfies $\varphi(\bar{x}, \bar{y}) \rightarrow \exists \bar{z} \psi(\bar{x}, \bar{z})$ if for every (\bar{a}, \bar{b}) such that I satisfies $\varphi(\bar{a}, \bar{b})$, there is a tuple \bar{c} such that J satisfies $\psi(\bar{a}, \bar{c})$.

J is a **solution** for I

Example: Data exchange problem

Previous example: $Book(x_1, x_2, y_1) \rightarrow \exists z_1 Writer(x_2, x_1, z_1)$

<i>Book</i>	<i>Title</i>	<i>AuthorName</i>	<i>Affiliation</i>
<i>I</i> :	Algebra	Hungerford	U. Washington
	Real Analysis	Royden	Stanford

Possible solutions:

Example: Data exchange problem

Previous example: $Book(x_1, x_2, y_1) \rightarrow \exists z_1 Writer(x_2, x_1, z_1)$

<i>Book</i>	<i>Title</i>	<i>AuthorName</i>	<i>Affiliation</i>
I	Algebra	Hungerford	U. Washington
	Real Analysis	Royden	Stanford

Possible solutions:

<i>Writer</i>	<i>Name</i>	<i>BookTitle</i>	<i>Year</i>
J_1	Hungerford	Algebra	1974
	Royden	Real Analysis	1988

Example: Data exchange problem

Previous example: $Book(x_1, x_2, y_1) \rightarrow \exists z_1 Writer(x_2, x_1, z_1)$

I	<table><thead><tr><th><i>Book</i></th><th><i>Title</i></th><th><i>AuthorName</i></th><th><i>Affiliation</i></th></tr></thead><tbody><tr><td></td><td>Algebra</td><td>Hungerford</td><td>U. Washington</td></tr><tr><td></td><td>Real Analysis</td><td>Royden</td><td>Stanford</td></tr></tbody></table>	<i>Book</i>	<i>Title</i>	<i>AuthorName</i>	<i>Affiliation</i>		Algebra	Hungerford	U. Washington		Real Analysis	Royden	Stanford
<i>Book</i>	<i>Title</i>	<i>AuthorName</i>	<i>Affiliation</i>										
	Algebra	Hungerford	U. Washington										
	Real Analysis	Royden	Stanford										

Possible solutions:

J_1	<table><thead><tr><th><i>Writer</i></th><th><i>Name</i></th><th><i>BookTitle</i></th><th><i>Year</i></th></tr></thead><tbody><tr><td></td><td>Hungerford</td><td>Algebra</td><td>1974</td></tr><tr><td></td><td>Royden</td><td>Real Analysis</td><td>1988</td></tr></tbody></table>	<i>Writer</i>	<i>Name</i>	<i>BookTitle</i>	<i>Year</i>		Hungerford	Algebra	1974		Royden	Real Analysis	1988
<i>Writer</i>	<i>Name</i>	<i>BookTitle</i>	<i>Year</i>										
	Hungerford	Algebra	1974										
	Royden	Real Analysis	1988										

J_2	<table><thead><tr><th><i>Writer</i></th><th><i>Name</i></th><th><i>BookTitle</i></th><th><i>Year</i></th></tr></thead><tbody><tr><td></td><td>Hungerford</td><td>Algebra</td><td>\perp_1</td></tr><tr><td></td><td>Royden</td><td>Real Analysis</td><td>\perp_2</td></tr></tbody></table>	<i>Writer</i>	<i>Name</i>	<i>BookTitle</i>	<i>Year</i>		Hungerford	Algebra	\perp_1		Royden	Real Analysis	\perp_2
<i>Writer</i>	<i>Name</i>	<i>BookTitle</i>	<i>Year</i>										
	Hungerford	Algebra	\perp_1										
	Royden	Real Analysis	\perp_2										

Outline

- ▶ The data exchange scenario
 - ▶ Query answering
- ▶ Queries with negation in data exchange
- ▶ $\text{DATALOG}^{\mathbf{C}(\neq)}$ programs
- ▶ Beyond union of conjunctive queries
 - ▶ Expressive power of $\text{DATALOG}^{\mathbf{C}(\neq)}$
 - ▶ New tractable classes of queries
- ▶ Concluding remarks

Query answering in data exchange

Given: Data exchange setting $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$, a query Q over \mathbf{T} and an instance I of \mathbf{S} .

- ▶ What does it mean to answer Q ?

Query answering in data exchange

Given: Data exchange setting $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$, a query Q over \mathbf{T} and an instance I of \mathbf{S} .

- ▶ What does it mean to answer Q ?

$$\text{CERTAIN}_{\mathcal{M}}(Q, I) = \bigcap_{J \text{ is a solution for } I} Q(J)$$

Example: Query answering in data exchange

Previous example: $Book(x_1, x_2, y_1) \rightarrow \exists z_1 Writer(x_2, x_1, z_1)$

<i>Book</i>	<i>Title</i>	<i>AuthorName</i>	<i>Affiliation</i>
<i>I</i> :	Algebra	Hungerford	U. Washington
	Real Analysis	Royden	Stanford

$$\text{CERTAIN}_{\mathcal{M}}(\exists y \exists z \textit{Writer}(x, y, z), I) = \{\text{Hungerford, Royden}\}$$

Computing certain answers

- ▶ A data exchange setting $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$ and a query Q are assumed to be fixed.
- ▶ Problem to solve:
 - Input : Instance I of \mathbf{S} and a tuple \bar{t} from I
 - Question : Is $\bar{t} \in \text{CERTAIN}_{\mathcal{M}}(Q, I)$?

Computing certain answers

- ▶ A data exchange setting $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$ and a query Q are assumed to be fixed.
- ▶ Problem to solve:
 - Input : Instance I of \mathbf{S} and a tuple \bar{t} from I
 - Question : Is $\bar{t} \in \text{CERTAIN}_{\mathcal{M}}(Q, I)$?

We are considering the **data complexity** of the problem.

Computing certain answers (cont'd)

How can $\text{CERTAIN}_{\mathcal{M}}(Q, I)$ be computed?

- ▶ Naïve algorithm does not work: infinitely many solutions

Computing certain answers (cont'd)

How can $\text{CERTAIN}_{\mathcal{M}}(Q, I)$ be computed?

- ▶ Naïve algorithm does not work: infinitely many solutions

Approach proposed in [FKMP03]:

1. Materialize a solution J for I such that:

$$\text{CERTAIN}_{\mathcal{M}}(Q, I) = Q(J)$$

2. Compute $Q(J)$

Computing certain answers (cont'd)

How can $\text{CERTAIN}_{\mathcal{M}}(Q, I)$ be computed?

- ▶ Naïve algorithm does not work: infinitely many solutions

Approach proposed in [FKMP03]:

1. Materialize a solution J for I such that:

$$\text{CERTAIN}_{\mathcal{M}}(Q, I) = Q(J)$$

2. Compute $Q(J)$

This works well for positive queries!

A solution to materialize: Canonical universal solution

Input: $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$ and an instance I of \mathbf{S}

Output: Canonical universal solution $\text{CAN}(I)$ for I

Algorithm:

for every $\varphi(\bar{x}, \bar{y}) \rightarrow \exists \bar{z} \psi(\bar{x}, \bar{z}) \in \Sigma$ do
 for every (\bar{a}, \bar{b}) such that I satisfies $\varphi(\bar{a}, \bar{b})$ do
 create a fresh tuple of null values \bar{n}
 insert $\psi(\bar{a}, \bar{n})$ into $\text{CAN}(I)$

Example: Canonical universal solution

Previous example: $Book(x_1, x_2, y_1) \rightarrow \exists z_1 Writer(x_2, x_1, z_1)$

<i>Book</i>	<i>Title</i>	<i>AuthorName</i>	<i>Affiliation</i>
<i>I</i> :	Algebra	Hungerford	U. Washington
	Real Analysis	Royden	Stanford

Example: Canonical universal solution

Previous example: $Book(x_1, x_2, y_1) \rightarrow \exists z_1 Writer(x_2, x_1, z_1)$

<i>Book</i>	<i>Title</i>	<i>AuthorName</i>	<i>Affiliation</i>
<i>I</i> :	Algebra	Hungerford	U. Washington
	Real Analysis	Royden	Stanford

We have that:

<i>Writer</i>	<i>Name</i>	<i>BookTitle</i>	<i>Year</i>
$CAN(I)$:	Hungerford	Algebra	\perp_1
	Royden	Real Analysis	\perp_2

Canonical universal solution: Computing certain answers

Canonical universal solution can be computed in **polynomial time**.

- ▶ Data complexity: Data exchange setting is fixed.

Canonical universal solution: Computing certain answers

Canonical universal solution can be computed in **polynomial time**.

- ▶ Data complexity: Data exchange setting is fixed.

Notation: $\mathbf{C}(a)$ holds if and only if a is a constant.

Canonical universal solution: Computing certain answers

Canonical universal solution can be computed in **polynomial time**.

- ▶ Data complexity: Data exchange setting is fixed.

Notation: $\mathbf{C}(a)$ holds if and only if a is a constant.

Theorem (FKMP03)

Let $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$, $Q(x_1, \dots, x_k)$ a union of conjunctive queries over \mathbf{T} and

$$Q^*(x_1, \dots, x_k) = \mathbf{C}(x_1) \wedge \dots \wedge \mathbf{C}(x_k) \wedge Q(x_1, \dots, x_k).$$

Then for every instance I of \mathbf{S} : $\text{CERTAIN}_{\mathcal{M}}(Q, I) = Q^*(\text{CAN}(I))$.

Why does the previous approach work?

Simple explanation: Closure under homomorphisms

Why does the previous approach work?

Simple explanation: Closure under homomorphisms

$h : \text{dom}(J_1) \rightarrow \text{dom}(J_2)$ is a *homomorphism* from J_1 to J_2 if:

- ▶ h preserves the relations: If $R(a_1, \dots, a_k)$ is in J_1 , then $R(h(a_1), \dots, h(a_k))$ is in J_2 .
- ▶ h is the identity on constants.

A solution J for I under \mathcal{M} is *universal* if:

- ▶ For every solution J' for I under \mathcal{M} , there exists a homomorphism from J to J' .

$\text{CAN}(I)$ is a universal solution for I

Why does the previous approach work? (cont'd)

Theorem (FKMP03)

Let $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$, $Q(x_1, \dots, x_k)$ a union of conjunctive queries over \mathbf{T} and

$$Q^*(x_1, \dots, x_k) = \mathbf{C}(x_1) \wedge \dots \wedge \mathbf{C}(x_k) \wedge Q(x_1, \dots, x_k).$$

Then for every instance I of \mathbf{S} and universal solution J for I under \mathcal{M} : $\text{CERTAIN}_{\mathcal{M}}(Q, I) = Q^*(J)$.

Why does the previous approach work? (cont'd)

Theorem (FKMP03)

Let $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$, $Q(x_1, \dots, x_k)$ a union of conjunctive queries over \mathbf{T} and

$$Q^*(x_1, \dots, x_k) = \mathbf{C}(x_1) \wedge \dots \wedge \mathbf{C}(x_k) \wedge Q(x_1, \dots, x_k).$$

Then for every instance I of \mathbf{S} and universal solution J for I under \mathcal{M} : $\text{CERTAIN}_{\mathcal{M}}(Q, I) = Q^*(J)$.

Proof: From the fact that Q^* is closed under homomorphisms

DATALOG as a query language for data exchange systems

The previous approach works for any language closed under homomorphisms.

- ▶ DATALOG queries can also be computed in polynomial time.

DATALOG as a query language for data exchange systems

The previous approach works for any language closed under homomorphisms.

- ▶ DATALOG queries can also be computed in polynomial time.

Unfortunately, both DATALOG and union of conjunctive queries keep us on the realm of positive.

Outline

- ▶ The data exchange scenario
 - ▶ Query answering
- ▶ Queries with negation in data exchange
- ▶ $\text{DATALOG}^{\mathbf{C}(\neq)}$ programs
- ▶ Beyond union of conjunctive queries
 - ▶ Expressive power of $\text{DATALOG}^{\mathbf{C}(\neq)}$
 - ▶ New tractable classes of queries
- ▶ Concluding remarks

Are queries with negation interesting in data exchange?

$$\mathcal{M} : \begin{array}{l} G(x, y) \rightarrow E(x, y) \\ S(x) \rightarrow P(x) \\ T(x) \rightarrow R(x) \end{array}$$

$$Q : \exists x \exists y \exists z (E(x, z) \wedge E(z, y) \wedge \neg E(x, y))$$

Are queries with negation interesting in data exchange?

$$\mathcal{M} : \begin{array}{l} G(x, y) \rightarrow E(x, y) \\ S(x) \rightarrow P(x) \\ T(x) \rightarrow R(x) \end{array}$$

$$Q : \exists x \exists y \exists z (E(x, z) \wedge E(z, y) \wedge \neg E(x, y))$$

$$I : G(a, b), G(b, c)$$

Are queries with negation interesting in data exchange?

$$\begin{array}{l} \mathcal{M} : \\ \quad G(x, y) \rightarrow E(x, y) \\ \quad S(x) \rightarrow P(x) \\ \quad T(x) \rightarrow R(x) \end{array}$$

$$Q : \exists x \exists y \exists z (E(x, z) \wedge E(z, y) \wedge \neg E(x, y))$$

$$I : G(a, b), G(b, c)$$

$$J_1 : E(a, b), E(b, c)$$

Are queries with negation interesting in data exchange?

$$\begin{array}{l} \mathcal{M} : \\ \quad G(x, y) \rightarrow E(x, y) \\ \quad S(x) \rightarrow P(x) \\ \quad T(x) \rightarrow R(x) \end{array}$$

$$Q : \exists x \exists y \exists z (E(x, z) \wedge E(z, y) \wedge \neg E(x, y))$$

$$\begin{array}{l} I : G(a, b), G(b, c) \\ J_1 : E(a, b), E(b, c) \\ J_2 : E(a, b), E(b, c), E(a, c) \end{array}$$

Are queries with negation interesting in data exchange?

$$\begin{aligned} \mathcal{M} : \quad & G(x, y) \rightarrow E(x, y) \\ & S(x) \rightarrow P(x) \\ & T(x) \rightarrow R(x) \end{aligned}$$

$$Q : \exists x \exists y \exists z (E(x, z) \wedge E(z, y) \wedge \neg E(x, y))$$

$$\begin{aligned} I : & G(a, b), G(b, c) \\ J_1 : & E(a, b), E(b, c) \\ J_2 : & E(a, b), E(b, c), E(a, c) \end{aligned}$$

- ▶ Q is always false in a solution where E represents the transitive closure of G

But if positive queries are also considered ...

$$\mathcal{M} : \begin{array}{l} G(x, y) \rightarrow E(x, y) \\ S(x) \rightarrow P(x) \\ T(x) \rightarrow R(x) \end{array}$$

$$Q : \begin{array}{l} \exists x \exists y (P(x) \wedge R(y) \wedge E(x, y)) \vee \\ \exists x \exists y \exists z (E(x, z) \wedge E(z, y) \wedge \neg E(x, y)) \end{array}$$

But if positive queries are also considered ...

$$\mathcal{M} : \begin{array}{l} G(x, y) \rightarrow E(x, y) \\ S(x) \rightarrow P(x) \\ T(x) \rightarrow R(x) \end{array}$$

$$Q : \exists x \exists y (P(x) \wedge R(y) \wedge E(x, y)) \vee \\ \exists x \exists y \exists z (E(x, z) \wedge E(z, y) \wedge \neg E(x, y))$$

$\text{CERTAIN}_{\mathcal{M}}(Q, I) = \text{true}$ iff there exist a, b such that:

- ▶ $P(a)$ and $R(b)$ hold in I
- ▶ (a, b) is in the transitive closure of G in I

Is there a *reasonable* query language with negation in data exchange?

We cannot directly add inequalities or negated relational atoms to DATALOG.

- ▶ Preservation under homomorphisms is lost
- ▶ Language becomes intractable, even for conjunctive queries with inequalities (AD98)

Is there a *reasonable* query language with negation in data exchange?

We cannot directly add inequalities or negated relational atoms to `DATALOG`.

- ▶ Preservation under homomorphisms is lost
- ▶ Language becomes intractable, even for conjunctive queries with inequalities (AD98)

Homomorphisms in data exchange are the identity on constants

- ▶ Inequalities witnessed by constants are preserved under homomorphisms

Our contributions

Query language that extends DATALOG with a form of negation

- ▶ As good as DATALOG for data exchange
- ▶ Can be used to find new tractable classes of queries

Outline

- ▶ The data exchange scenario
 - ▶ Query answering
- ▶ Queries with negation in data exchange
- ▶ $\text{DATALOG}^{\mathbf{C}(\neq)}$ programs
- ▶ Beyond union of conjunctive queries
 - ▶ Expressive power of $\text{DATALOG}^{\mathbf{C}(\neq)}$
 - ▶ New tractable classes of queries
- ▶ Concluding remarks

DATALOG^{C(≠)} programs

Definition

A *constant-inequality Datalog rule* is a rule of the form:

$$S(\bar{x}) \leftarrow S_1(\bar{x}_1), \dots, S_\ell(\bar{x}_\ell), \mathbf{C}(y_1), \dots, \mathbf{C}(y_m), u_1 \neq v_1, \dots, u_n \neq v_n$$

where

- ▶ S, S_1, \dots, S_ℓ are predicate symbols,
- ▶ every variable in \bar{x} is mentioned in some tuple \bar{x}_i ,
- ▶ every variable y_j is mentioned in some tuple \bar{x}_i , and
- ▶ every variable u_j , and every variable v_j , is equal to some variable y_i .

DATALOG^{C(≠)} programs (cont'd)

A *constant-inequality Datalog program* (DATALOG^{C(≠)} program) is a finite set of constant-inequality Datalog rules.

DATALOG^{C(≠)} programs (cont'd)

A *constant-inequality Datalog program* (DATALOG^{C(≠)} program) is a finite set of constant-inequality Datalog rules.

Example:

$$S(x, y) \leftarrow E(x, y)$$

$$S(x, y) \leftarrow S(x, u), S(u, v), S(v, y), \mathbf{C}(x), \mathbf{C}(u), \mathbf{C}(v), u \neq v$$

DATALOG^{C(≠)} programs (cont'd)

A *constant-inequality Datalog program* (DATALOG^{C(≠)} program) is a finite set of constant-inequality Datalog rules.

Example:

$$S(x, y) \leftarrow E(x, y)$$

$$S(x, y) \leftarrow S(x, u), S(u, v), S(v, y), \mathbf{C(x)}, \mathbf{C(u)}, \mathbf{C(v)}, u \neq v$$

DATALOG^{C(≠)} programs (cont'd)

A *constant-inequality Datalog program* (DATALOG^{C(≠)} program) is a finite set of constant-inequality Datalog rules.

Example:

$$S(x, y) \leftarrow E(x, y)$$

$$S(x, y) \leftarrow S(x, u), S(u, v), S(v, y), \mathbf{C}(x), \mathbf{C}(u), \mathbf{C}(v), u \neq v$$

$\text{DATALOG}^{\mathbf{C}(\neq)}$ programs can be evaluated efficiently

$\text{DATALOG}^{\mathbf{C}(\neq)}$ programs are preserved under homomorphisms

$\text{DATALOG}^{\mathbf{C}(\neq)}$ programs can be evaluated efficiently

$\text{DATALOG}^{\mathbf{C}(\neq)}$ programs are preserved under homomorphisms

Proposition

Certain answers of $\text{DATALOG}^{\mathbf{C}(\neq)}$ programs can be computed by evaluating the programs over the canonical universal solution.

DATALOG^{C(≠)} programs can be evaluated efficiently

DATALOG^{C(≠)} programs are preserved under homomorphisms

Proposition

Certain answers of DATALOG^{C(≠)} programs can be computed by evaluating the programs over the canonical universal solution.

Theorem

Computing the certain answers of a DATALOG^{C(≠)} program takes polynomial time (data complexity)

Outline

- ▶ The data exchange scenario
 - ▶ Query answering
- ▶ Queries with negation in data exchange
- ▶ $\text{DATALOG}^{\mathbf{C}(\neq)}$ programs
- ▶ Beyond union of conjunctive queries
 - ▶ Expressive power of $\text{DATALOG}^{\mathbf{C}(\neq)}$
 - ▶ New tractable classes of queries
- ▶ Concluding remarks

DATALOG^{C(≠)} programs can express queries with negation

Theorem (ABR09)

For every union of conjunctive queries Q with at most one

- ▶ negated relational atom or*
- ▶ inequality*

per disjunct, there exists a DATALOG^{C(≠)} program Π such that

Q and Π are equivalent in the data exchange scenario.

DATALOG^{C(≠)} programs can express queries with negation

Theorem (ABR09)

For every union of conjunctive queries Q with at most one

- ▶ negated relational atom or*
- ▶ inequality*

per disjunct, there exists a DATALOG^{C(≠)} program Π such that

for every data exchange setting $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$ and instance I of \mathbf{S} : $\text{CERTAIN}_{\mathcal{M}}(Q, I) = \text{CERTAIN}_{\mathcal{M}}(\Pi, I)$.

DATALOG^{C(≠)} programs can express queries with negation

Theorem (ABR09)

For every union of conjunctive queries Q with at most one

- ▶ negated relational atom or*
- ▶ inequality*

per disjunct, there exists a DATALOG^{C(≠)} program Π such that

for every data exchange setting $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$ and instance I of \mathbf{S} : $\text{CERTAIN}_{\mathcal{M}}(Q, I) = \text{CERTAIN}_{\mathcal{M}}(\Pi, I)$.

- ▶ Certain answers for this class of queries can be computed in polynomial time.
 - ▶ Π can be computed from Q in polynomial time.
- ▶ Result for inequalities was proved in [FKMP03].

Example: Expressing negation in DATALOG^{C(≠)}

$$Q : \exists x \exists y (E(x, y) \wedge x \neq y) \vee \\ \exists x \exists y \exists z (E(x, y) \wedge E(y, z) \wedge \neg E(x, z))$$

Example: Expressing negation in DATALOG^{C(≠)}

$$Q : \exists x \exists y (E(x, y) \wedge x \neq y) \vee \\ \exists x \exists y \exists z (E(x, y) \wedge E(y, z) \wedge \neg E(x, z))$$

$$\begin{aligned} \text{dom}(x) &\leftarrow E(x, z) && \text{Collect the domain} \\ \text{dom}(x) &\leftarrow E(z, x) \end{aligned}$$

Example: Expressing negation in DATALOG^{C(≠)}

$$Q : \exists x \exists y (E(x, y) \wedge x \neq y) \vee \\ \exists x \exists y \exists z (E(x, y) \wedge E(y, z) \wedge \neg E(x, z))$$

$dom(x) \leftarrow E(x, z)$ Formalize equality

$dom(x) \leftarrow E(z, x)$

$EQ(x, x) \leftarrow dom(x)$

$EQ(x, y) \leftarrow EQ(y, x)$

$EQ(x, y) \leftarrow EQ(x, z), EQ(z, y)$

Example: Expressing negation in DATALOG^{C(≠)}

$$Q : \exists x \exists y (E(x, y) \wedge x \neq y) \vee \\ \exists x \exists y \exists z (E(x, y) \wedge E(y, z) \wedge \neg E(x, z))$$

$$\begin{array}{lll} \text{dom}(x) & \leftarrow & E(x, z) \\ \text{dom}(x) & \leftarrow & E(z, x) \\ \text{EQ}(x, x) & \leftarrow & \text{dom}(x) \\ \text{EQ}(x, y) & \leftarrow & \text{EQ}(y, x) \\ \text{EQ}(x, y) & \leftarrow & \text{EQ}(x, z), \text{EQ}(z, y) \\ U(x, y) & \leftarrow & E(x, y) \end{array} \quad \text{Copy } E \text{ into } U$$

Example: Expressing negation in DATALOG^{C(≠)}

$$Q : \exists x \exists y (E(x, y) \wedge x \neq y) \vee \\ \exists x \exists y \exists z (E(x, y) \wedge E(y, z) \wedge \neg E(x, z))$$

$dom(x) \leftarrow E(x, z)$ Replace equals in U

$dom(x) \leftarrow E(z, x)$

$EQ(x, x) \leftarrow dom(x)$

$EQ(x, y) \leftarrow EQ(y, x)$

$EQ(x, y) \leftarrow EQ(x, z), EQ(z, y)$

$U(x, y) \leftarrow E(x, y)$

$U(x, y) \leftarrow U(u, v), EQ(u, x), EQ(v, y)$

Example: Expressing negation in DATALOG^{C(≠)}

$$Q : \exists x \exists y (E(x, y) \wedge x \neq y) \vee \\ \exists x \exists y \exists z (E(x, y) \wedge E(y, z) \wedge \neg E(x, z))$$

$dom(x) \leftarrow E(x, z)$ Simulate 2nd disjunct of Q

$dom(x) \leftarrow E(z, x)$

$EQ(x, x) \leftarrow dom(x)$

$EQ(x, y) \leftarrow EQ(y, x)$

$EQ(x, y) \leftarrow EQ(x, z), EQ(z, y)$

$U(x, y) \leftarrow E(x, y)$

$U(x, y) \leftarrow U(u, v), EQ(u, x), EQ(v, y)$

$U(x, y) \leftarrow U(x, z), U(z, y)$

Example: Expressing negation in DATALOG^{C(≠)}

$$Q : \exists x \exists y (E(x, y) \wedge x \neq y) \vee \\ \exists x \exists y \exists z (E(x, y) \wedge E(y, z) \wedge \neg E(x, z))$$

dom(*x*) ← *E*(*x*, *z*) Simulate 1st disjunct of *Q*
dom(*x*) ← *E*(*z*, *x*)
EQ(*x*, *x*) ← *dom*(*x*)
EQ(*x*, *y*) ← *EQ*(*y*, *x*)
EQ(*x*, *y*) ← *EQ*(*x*, *z*), *EQ*(*z*, *y*)
U(*x*, *y*) ← *E*(*x*, *y*)
U(*x*, *y*) ← *U*(*u*, *v*), *EQ*(*u*, *x*), *EQ*(*v*, *y*)
U(*x*, *y*) ← *U*(*x*, *z*), *U*(*z*, *y*)
EQ(*x*, *y*) ← *U*(*x*, *y*)

Example: Expressing negation in DATALOG^{C(≠)}

$Q : \exists x \exists y (E(x, y) \wedge x \neq y) \vee$
 $\exists x \exists y \exists z (E(x, y) \wedge E(y, z) \wedge \neg E(x, z))$

$dom(x) \leftarrow E(x, z)$ Answer Q
 $dom(x) \leftarrow E(z, x)$
 $EQ(x, x) \leftarrow dom(x)$
 $EQ(x, y) \leftarrow EQ(y, x)$
 $EQ(x, y) \leftarrow EQ(x, z), EQ(z, y)$
 $U(x, y) \leftarrow E(x, y)$
 $U(x, y) \leftarrow U(u, v), EQ(u, x), EQ(v, y)$
 $U(x, y) \leftarrow U(x, z), U(z, y)$
 $EQ(x, y) \leftarrow U(x, y)$
TRUE $\leftarrow EQ(x, y), C(x), C(y), x \neq y$

Outline

- ▶ The data exchange scenario
 - ▶ Query answering
- ▶ Queries with negation in data exchange
- ▶ $\text{DATALOG}^{\mathbf{C}(\neq)}$ programs
- ▶ Beyond union of conjunctive queries
 - ▶ Expressive power of $\text{DATALOG}^{\mathbf{C}(\neq)}$
 - ▶ **New tractable classes of queries**
- ▶ Concluding remarks

Certain answers for conjunctive queries with two inequalities

Theorem (M05)

The certain answers problem is coNP -complete for the class of conjunctive queries with two inequalities (data complexity).

Certain answers for conjunctive queries with two inequalities

Theorem (M05)

The certain answers problem is coNP -complete for the class of conjunctive queries with two inequalities (data complexity).

An interesting tractable fragment of 2-CQ^{\neq} was identified by considering $\text{DATALOG}^{\text{C}(\neq)}$ programs.

Two restrictions on queries are needed

Just the intuition:

Two restrictions on queries are needed

Just the intuition:

- ▶ **Constant Joins:** Null values do not witness a join

Two restrictions on queries are needed

Just the intuition:

- ▶ **Constant Joins:** Null values do not witness a join
- ▶ **Almost constant inequalities:** Every inequality is not witnessed just by null values

A tractable fragment of 2-UCQ \neq

Theorem (ABR09)

For every 2-UCQ \neq Q with:

- ▶ *constant joins and*
- ▶ *almost constant inequalities,*

there exists a DATALOG $^{\mathbf{C}(\neq)}$ program Π such that:

Q and Π are equivalent in the data exchange scenario.

A tractable fragment of 2-UCQ \neq

Theorem (ABR09)

For every 2-UCQ \neq Q with:

- ▶ constant joins and
- ▶ almost constant inequalities,

there exists a DATALOG $^{\mathbf{C}(\neq)}$ program Π such that:

for every data exchange setting $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$ and instance I of \mathbf{S} : $\text{CERTAIN}_{\mathcal{M}}(Q, I) = \text{CERTAIN}_{\mathcal{M}}(\Pi, I)$.

A tractable fragment of 2-UCQ \neq

Theorem (ABR09)

For every 2-UCQ \neq Q with:

- ▶ constant joins and
- ▶ almost constant inequalities,

there exists a DATALOG $^{\mathbf{C}(\neq)}$ program Π such that:

for every data exchange setting $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$ and instance I of \mathbf{S} : $\text{CERTAIN}_{\mathcal{M}}(Q, I) = \text{CERTAIN}_{\mathcal{M}}(\Pi, I)$.

Certain answers to this class of queries can be computed in polynomial time

- ▶ Π can be computed from Q in polynomial time.

A tractable fragment of 2-UCQ \neq (cont'd)

Theorem (ABR09)

Removing any one of the conditions in the previous theorem yields to CONP -completeness.

A tractable fragment of 2-UCQ[≠] (cont'd)

Theorem (ABR09)

Removing any one of the conditions in the previous theorem yields to CONP-completeness.

This result is stronger than the result in [M05].

- ▶ Reduction in [M05] does not impose these restrictions.

Outline

- ▶ The data exchange scenario
 - ▶ Query answering
- ▶ Queries with negation in data exchange
- ▶ $\text{DATALOG}^{\mathbf{C}(\neq)}$ programs
- ▶ Beyond union of conjunctive queries
 - ▶ Expressive power of $\text{DATALOG}^{\mathbf{C}(\neq)}$
 - ▶ New tractable classes of queries
- ▶ Concluding remarks

We propose $\text{DATALOG}^{\mathbf{C}(\neq)}$ as a query language for data exchange systems

- ▶ Certain answers to a $\text{DATALOG}^{\mathbf{C}(\neq)}$ program can be computed in polynomial time (data complexity).
- ▶ $\text{DATALOG}^{\mathbf{C}(\neq)}$ is equipped with a form of negation.
 - ▶ Union of conjunctive queries with one negated atom per disjunct are expressible in $\text{DATALOG}^{\mathbf{C}(\neq)}$.
- ▶ $\text{DATALOG}^{\mathbf{C}(\neq)}$ can be used to find new tractable classes of queries with negation.
 - ▶ In this talk: A fragment of 2-UCQ^{\neq}

Thank you!

Two restrictions on queries are needed

Two restrictions on queries are needed

Constant Joins

Null values do not witness a join

$$\mathcal{M} : \begin{array}{l} P(u, v) \rightarrow T(u, v) \\ Q(u, v) \rightarrow \exists w U(u, w) \end{array}$$

$$Q_1 : \exists x \exists y \exists z (T(x, y) \wedge U(x, z))$$

$$Q_2 : \exists x \exists y \exists z (U(x, z) \wedge U(y, z))$$

Two restrictions on queries are needed

Constant Joins

Null values do not witness a join

$$\mathcal{M} : \begin{array}{l} P(u, v) \rightarrow T(u, v) \\ Q(u, v) \rightarrow \exists w U(u, w) \end{array}$$

$$Q_1 : \exists x \exists y \exists z (T(x, y) \wedge U(x, z))$$

YES

$$Q_2 : \exists x \exists y \exists z (U(x, z) \wedge U(y, z))$$

Two restrictions on queries are needed

Constant Joins

Null values do not witness a join

$$\mathcal{M} : \begin{array}{l} P(u, v) \rightarrow T(u, v) \\ Q(u, v) \rightarrow \exists w U(u, w) \end{array}$$

$$Q_1 : \exists x \exists y \exists z (T(x, y) \wedge U(x, z)) \quad \text{YES}$$

$$Q_2 : \exists x \exists y \exists z (U(x, z) \wedge U(y, z)) \quad \text{NO}$$

Two restrictions on queries are needed

Almost constant inequalities

Every inequality is not witnessed just by null values

$$\mathcal{M} : \begin{array}{l} P(u, v) \rightarrow T(u, v) \\ Q(u, v) \rightarrow \exists w U(u, w) \end{array}$$

$$Q_1 : \exists x \exists y \exists z (U(x, y) \wedge U(x, z) \wedge x \neq z)$$

$$Q_2 : \exists x \exists y \exists z (U(x, y) \wedge U(x, z) \wedge y \neq z)$$

Two restrictions on queries are needed

Almost constant inequalities

Every inequality is not witnessed just by null values

$$\mathcal{M} : \begin{array}{l} P(u, v) \rightarrow T(u, v) \\ Q(u, v) \rightarrow \exists w U(u, w) \end{array}$$

$$Q_1 : \exists x \exists y \exists z (U(x, y) \wedge U(x, z) \wedge x \neq z)$$

YES

$$Q_2 : \exists x \exists y \exists z (U(x, y) \wedge U(x, z) \wedge y \neq z)$$

Two restrictions on queries are needed

Almost constant inequalities

Every inequality is not witnessed just by null values

$$\mathcal{M} : \begin{array}{l} P(u, v) \rightarrow T(u, v) \\ Q(u, v) \rightarrow \exists w U(u, w) \end{array}$$

$$Q_1 : \exists x \exists y \exists z (U(x, y) \wedge U(x, z) \wedge x \neq z) \quad \text{YES}$$

$$Q_2 : \exists x \exists y \exists z (U(x, y) \wedge U(x, z) \wedge y \neq z) \quad \text{NO}$$