

Computing Range Consistent Answers to Aggregation Queries via Rewriting ¹

Aziz Amezian El Khalfioui and Jef Wijsen

University of Mons, Belgium

November 22, 2024

¹This paper will also be presented at PODS 2025 [AEKW24].

Table of Contents

- 1 Numerical queries
- 2 Inconsistent Databases
- 3 Range Consistent Answers
- 4 Main Result

Numerical queries

We consider numerical queries that take the following form

$$\text{AGG}(r) \leftarrow q(\vec{u})$$

where

- AGG is an aggregate symbol (like MAX , MIN , SUM , AVG , COUNT);
- $q(\vec{u})$ is a self-join-free conjunction of atoms with variables \vec{u} ; and
- r is either a numeric variable occurring in \vec{u} or a constant rational number.

Definition

Let \mathbf{db} be a database instance. An embedding of $q(\vec{u})$ in \mathbf{db} is a total mapping θ from \vec{u} to the set of constants appearing in \mathbf{db} such that $\theta(q) \subseteq \mathbf{db}$.

Semantics of Numerical Queries

Example

Consider the conjunction $Dealers(\underline{x}, y) \wedge Stock(\underline{z}, y, r)$, and the following database instance **db**:

<i>Dealers</i>	<u>Name</u>	<u>Town</u>	<i>Stock</i>	<u>Product</u>	<u>Town</u>	<u>Qty</u>
	Smith	Boston		Tesla X	Boston	35
	Smith	New York		Tesla Y	Boston	35
				Tesla Y	New York	90

We have 3 different embeddings in **db**:

- one for Tesla X sold in Boston;
- one for Tesla Y sold in Boston; and
- one for Tesla Y sold in New York.

Example

Consider the numerical query

$$g() = \text{SUM}(r) \leftarrow \text{Dealers}(\underline{x}, y) \wedge \text{Stock}(\underline{z}, y, r)$$

and the following database instance **db**:

<i>Dealers</i>	<u>Name</u> <u>Town</u>		<i>Stock</i>	<u>Product</u>	<u>Town</u>	<u>Qty</u>
	<u>Name</u>	<u>Town</u>		<u>Product</u>	<u>Town</u>	<u>Qty</u>
	Smith	Boston		Tesla X	Boston	35
	Smith	New York		Tesla Y	Boston	35
				Tesla Y	New York	90

- $g()$ returns the sum over the multiset containing $\theta(r)$ for each embedding θ in **db**.
- In this example, $g()$ returns $\text{SUM}(\{\{35, 35, 90\}\}) = 160$.

Table of Contents

- 1 Numerical queries
- 2 Inconsistent Databases**
- 3 Range Consistent Answers
- 4 Main Result

Inconsistent Databases

Primary Key Violations

We consider database instances that may violate primary key constraints.

Example

Consider the following database instance:

<i>Dealers</i>	<u><i>Name</i></u>	<u><i>Town</i></u>	<i>Stock</i>	<u><i>Product</i></u>	<u><i>Town</i></u>	<u><i>Qty</i></u>
	Smith	Boston		Tesla X	Boston	35
	Smith	New York		Tesla Y	Boston	35
				Tesla Y	New York	90

The primary key $Name \rightarrow Town$ is violated.

Inconsistent Databases

Definition

Let **db** be a database instance. A **repair** is an inclusion-maximal subset that satisfies all primary keys.

Example

Consider the following database instance:

<i>Dealers</i>	<u>Name</u>	<u>Town</u>	<i>Stock</i>	<u>Product</u>	<u>Town</u>	<u>Qty</u>
	Smith	Boston		Tesla X	Boston	35
	Smith	New York		Tesla Y	Boston	35
				Tesla Y	New York	90

We have two repairs:

- one where the dealer Smith works in Boston; and
- one where the dealer Smith works in New York.

Table of Contents

- 1 Numerical queries
- 2 Inconsistent Databases
- 3 Range Consistent Answers**
- 4 Main Result

Range Consistent Answers

We employ the range semantics as presented by Arenas et al. [ABC01], which provides the greatest lower bound (glb) and the least upper bound (lub) of query answers across all repairs.

Example

Consider the numerical query $SUM(r) \leftarrow Dealers(\underline{x}, y) \wedge Stock(\underline{z}, y, r)$ and the following database instance:

<i>Dealers</i>	<i>Dealers</i>		<i>Stock</i>	<i>Product</i>	<i>Town</i>	<i>Qty</i>
	<u>Name</u>	<u>Town</u>				
	Smith	Boston		Tesla X	Boston	35
	Smith	New York		Tesla Y	Boston	35
				Tesla Y	New York	90

- The numerical query returns $SUM(\{\{35, 35\}\}) = 70$ on one repair, and $SUM(\{\{90\}\}) = 90$ on the other repair.
- Thus, the glb is 70, the lub is 90, and the range consistent answer is $[70, 90]$.

Definition

Let $g()$ be a numerical query. The function problems $\text{GLB-CQA}(g())$ and $\text{LUB-CQA}(g())$ take a database instance \mathbf{db} as input, and return, respectively, the glb and the lub of the set that contains each number returned by $g()$ on some repair.

In this work, we focus on $\text{GLB-CQA}(g())$.

Table of Contents

- 1 Numerical queries
- 2 Inconsistent Databases
- 3 Range Consistent Answers
- 4 Main Result**

Main Result

Definition

An aggregate operator is **associative** if for all non-empty multisets X and Y such that $X \neq \emptyset$, we have $\mathcal{F}(X \uplus Y) = \mathcal{F}(\{\mathcal{F}(X)\} \uplus Y)$, where \uplus denotes union of multisets.

Definition

An aggregate operator \mathcal{F} is **monotone** if for all $m > 0$ and every (possibly empty) multiset Y , we have $\mathcal{F}(\{x_1, \dots, x_m\}) \leq \mathcal{F}(\{x'_1, \dots, x'_m\} \uplus Y)$ whenever $x_i \leq x'_i$ for every i .

Example

MAX and SUM (over non-negative values) are monotone and associative.
MIN is associative but not monotone since $\text{MIN}(\{2, 3\}) \not\leq \text{MIN}(\{2, 3, 1\})$.

The logic AGGR[FOL] extends FOL with aggregation operators.

Theorem

The following decision problem is decidable in quadratic time (in the size of the input):

*Given a numerical query $g() := \text{AGG}(r) \leftarrow q(\vec{u})$ such that the aggregate operator AGG is **monotone** and **associative**, is GLB-CQA($g()$) expressible in AGGR[FOL]?*

Moreover, if the answer is “yes,” it is possible to effectively construct, also in quadratic time, a formula in AGGR[FOL] that solves GLB-CQA($g()$).

- Extend our results to aggregate operators that lack monotonicity, associativity, or both.
- Shift from expressibility in $\text{AGGR}[\text{FOL}]$ to computability in P .

Thanks!

References I



Marcelo Arenas, Leopoldo E. Bertossi, and Jan Chomicki.

Scalar aggregation in fd-inconsistent databases.

In *ICDT*, volume 1973 of *Lecture Notes in Computer Science*, pages 39–53. Springer, 2001.



Aziz Amezian El Khalfioui and Jef Wijsen.

Computing range consistent answers to aggregation queries via rewriting.

Proc. ACM Manag. Data, 2(5), November 2024.