

Simple Aggregations in Algebraic Databases

Patrick Schultz, David I. Spivak, Ryan Wisnesky

Abstract

This document describes an extension to the Algebraic Databases formalism (Schultz & Wisnesky, 2017) that allows for simple aggregations in uber-flower queries.

1 Extending Multi-sorted Equational Logic

In this paper we define a syntax, semantics, and proof system that extends multi-sorted equational logic. This system is used in the AQL tool to implement simple aggregations in uber-flower queries.

1.1 Syntax

A *signature* Sig consists of:

1. A set *Sorts* whose elements are called *sorts*,
2. A set *Symbols* of pairs $(f, s_1 \times \dots \times s_k \rightarrow s)$ with $s_1, \dots, s_k, s \in \text{Sorts}$ and no f occurring in two distinct pairs. We write $f : X$ instead of $(f, X) \in \text{Symbols}$. When $k = 0$, we may call f a *constant symbol* and write $f : s$ instead of $f : \rightarrow s$. Otherwise, we may call f a *function symbol*.

We assume we have some countably infinite set $\{v_1, v_2, \dots\}$, whose elements we call *variables* and which are assumed to be distinct from any sort or symbol we ever consider. A *context* Γ is defined as a finite set of variable-sort pairs, with no variable given more than one sort:

$$\Gamma := \{v_1 : s_1, \dots, v_k : s_k\}$$

We inductively define the set $\text{Terms}^s(\text{Sig}, \Gamma)$ of *terms* of sort s over signature Sig and context Γ as:

1. $x \in \text{Terms}^s(\text{Sig}, \Gamma)$, if $x : s \in \Gamma$,
2. $f(t_1, \dots, t_k) \in \text{Terms}^s(\text{Sig}, \Gamma)$, if $f : s_1 \times \dots \times s_k \rightarrow s$ and $t_i \in \text{Terms}^{s_i}(\text{Sig}, \Gamma)$ for $i = 1, \dots, k$. When $k = 0$, we may write f for $f()$.
3. $(f, o)\{\text{for } \Gamma' \text{ where } t_1 = t'_1, \dots, t_k = t'_k \text{ return } e\} \in \text{Terms}^s(\text{Sig}, \Gamma)$ when Γ' is a context such that $\Gamma \cap \Gamma' = \emptyset$ and $f : s \times s \rightarrow s$ and $o : s$ and $e \in \text{Terms}^s(\text{Sig}, \Gamma \cup \Gamma')$ and for every $1 \leq i \leq k$ there exists a sort s_i such that $t_i, t'_i \in \text{Terms}^{s_i}(\text{Sig}, \Gamma \cup \Gamma')$.

Note that the monoid comprehensions (for Γ' -terms) are binding constructs: the variables in Γ' are considered bound. (Capture-avoiding) substitution of a term t for a variable v in a term e is written as $e[v \mapsto t]$ and is defined as usual.

2

An *equation* over Sig is a formula $\forall\Gamma. t_1 = t_2 : s$ with $t_1, t_2 \in Terms^s(Sig, \Gamma)$; we will omit the $: s$ when doing so will not lead to confusion. A *theory* is a pair of a signature and a set of equations over that signature. Associated with a theory Th is a binary relation between terms, called *provable equality*. We write $Th \vdash \forall\Gamma. t = t' : s$ to indicate that the theory Th proves that terms $t, t' \in Terms^s(Sig, \Gamma)$ are equal according to the usual rules of multi-sorted equational logic extended with five monoid comprehension laws. The equational logic rules are

$$\frac{t \in Terms^s(Sig, \Gamma)}{Th \vdash \forall\Gamma. t = t : s} \quad \frac{Th \vdash \forall\Gamma. t = t' : s}{Th \vdash \forall\Gamma. t' = t : s} \quad \frac{Th \vdash \forall\Gamma. t = t' : s \quad Th \vdash \forall\Gamma. t' = t'' : s}{Th \vdash \forall\Gamma. t = t'' : s}$$

$$\frac{Th \vdash \forall\Gamma. t = t' : s \quad v \notin \Gamma}{Th \vdash \forall\Gamma, v : s'. t = t' : s} \quad \frac{Th \vdash \forall\Gamma, v : s. t = t' : s' \quad Th \vdash \forall\Gamma. e = e' : s}{Th \vdash \forall\Gamma. t[v \mapsto e] = t'[v \mapsto e'] : s'}$$

and the monoid rules are

$$\frac{}{Th \vdash \forall\Gamma. (f, o)\{\text{for } - \text{ where } - \text{ return } e\} = e}$$

$$\frac{}{Th \vdash \forall\Gamma. (f, o)\{\text{for } \Gamma' \text{ where } \phi \text{ return } o\} = o}$$

$$\frac{Th \vdash \forall\Gamma. (f, o)\{\text{for } \Gamma' \text{ where } \phi \text{ return } f(e, e')\} = f((f, o)\{\text{for } \Gamma' \text{ where } \phi \text{ return } e\}, (f, o)\{\text{for } \Gamma' \text{ where } \phi \text{ return } e'\})}{Th \vdash \forall\Gamma. (f, o)\{\text{for } \Gamma' \text{ where } \phi \text{ return } (f, o)\{\text{for } \Gamma'' \text{ where } \phi' \text{ return } e\}\} = (f, o)\{\text{for } \Gamma' \cup \Gamma'' \text{ where } \phi \cup \phi' \text{ return } e\}}$$

$$\frac{Th \vdash h(o) = o' \quad Th \vdash \forall xy. h(f(x, y)) = f'(h(x), h(y))}{Th \vdash \forall\Gamma. h((f, o)\{\text{for } \Gamma' \text{ where } \phi \text{ return } e\}) = (f', o')\{\text{for } \Gamma' \text{ where } \phi \text{ return } h(e)\}}$$

We say that a theory Th is *OK for aggregation* when, for every term in Th of the form

$$(f, o)\{\text{for } \Gamma \text{ where } t \text{ return } e\}$$

we have

$$Th \vdash' \forall xyz. f(x, f(y, z)) = f(f(x, y), z) \quad Th \vdash' \forall x. f(x, o) = x = f(o, x) \quad Th \vdash' \forall xy. f(x, y) = f(y, x)$$

where $Th \vdash' \phi$ indicates that Th entails ϕ under the rules of multi-sorted equational logic without use of the monoid laws. We only consider theories that are OK for aggregation.

1.2 Semantics

An *algebra* A over a signature Sig consists of:

- a set $A(s)$ for each sort s ; the elements of $A(s)$ are called *carriers*, and

- a function $A(f) : A(s_1) \times \dots \times A(s_k) \rightarrow A(s)$ for each symbol $f : s_1 \times \dots \times s_k \rightarrow s$.

Let $\Gamma := \{x_1 : s_1, \dots, x_n : s_n\}$ be a context. An A -environment η for Γ associates each variable x_i with an element of $A(s_i)$. Write $\llbracket \Gamma \rrbracket$ to indicate the set of all A -environments for Γ .

The meaning of a term t in $Terms(Sig, \Gamma)$ relative to A -environment η for Γ is written $\llbracket t \rrbracket \eta$ and recursively defined as:

$$\llbracket x \rrbracket \eta = \eta(x) \quad \llbracket f(t_1, \dots, t_n) \rrbracket \eta = A(f)(\llbracket t_1 \rrbracket \eta, \dots, \llbracket t_n \rrbracket \eta)$$

To extend the above definition to aggregations

$$\llbracket (f, o)\{\text{for } \Gamma' \text{ where } t_1 = t'_1, \dots \text{ return } e\} \rrbracket \eta$$

first define the set

$$\zeta := \{\eta' \mid \eta' \in \llbracket \Gamma' \rrbracket, \llbracket t_1 \rrbracket \eta \cup \eta' = \llbracket t'_1 \rrbracket \eta \cup \eta', \dots\}$$

and the meaning of the agg term is

- $A(o)$ if ζ is empty,
- $\llbracket e \rrbracket \eta'_1 \cup \eta$ if $\zeta = \{\eta'_1\}$,
- $A(f)(\llbracket e \rrbracket \eta'_1 \cup \eta, \llbracket e \rrbracket \eta'_2 \cup \eta)$ if $\zeta = \{\eta'_1, \eta'_2\}$,
- etc

References

Patrick Shultz & Ryan Wisnesky *Algebraic Data Integration*.

7 March 2017