

Initial Value Computation

August 4, 2011

1 Definitions

In this draft we want to explain the notion of domains in DBToaster calculus [1]. DBToaster uses query language AGCA (which stands for AGgregation CAlculus). AGCA expressions are built from constants, variables, relational atoms, aggregate sums (Sum), conditions, and variable assignments (\leftarrow) using “+”, “”, and “.”. The abstract syntax can be given by the EBNF:

$$q ::= q \cdot q \mid q + q \mid v \leftarrow q \mid v_1 \theta v_2 \mid R(\vec{y}) \mid c \mid v \mid M[\vec{x}][\vec{y}] \quad (1)$$

The above definition can express a vast majority of all SQL statements. Here v denotes variables, \vec{x}, \vec{y} tuples of variables, R relation names, c constants, and θ denotes comparison operations ($=, \neq, >, \geq, <, \text{ and } \leq$). “+” represents unions and “.” is a DBToaster equivalent for joins. Assignment operator (\leftarrow) takes an query and assigns its result to a variable (v). A map $M[\vec{x}][\vec{y}]$ is a subquery with some input (\vec{x}) and output (\vec{y}) variables. It can be seen as a nested query which for variables \vec{x} produces the output \vec{y} , it is not defined in [1] but we added here for the purpose of this work.

The domain of a variable is the set of values that it can take. The domain of all the variables in a query expression can easily be computed recursively if some rules are respected. We will use through out the entire paper the notation of $\text{dom}_{\vec{x}}(q)$ for the domain of a set of variables, where q is the given query and \vec{x} is a vector representing the variables (not necessarily present in the expression q). We start with the definition of $\text{dom}_{\vec{x}}(R(\vec{y}))$:

$$\text{dom}_{\vec{x}}(R(\vec{y})) = \left\{ \vec{c} \mid \left(\prod_{i: x_i \in \vec{y}} (x_i \leftarrow c_i) \cdot R(\vec{y}) \right) \neq 0 \right\} \quad (2)$$

which in relation algebra means

$$\text{dom}_{\vec{x}}(R(\vec{y})) = \left\{ \vec{c} \mid \sigma_{i \in (\vec{x} \cap \vec{y}) : x_i = c_i} (\pi_{i : x_i \in \vec{y}} R(\vec{y})) \neq \text{NULL} \right\} \quad (3)$$

where $\vec{x} = \langle x_1, x_2, x_3, \dots, x_i, \dots \rangle$, $\vec{c} = \langle c_1, c_2, c_3, \dots, c_i, \dots \rangle$ in both of the definitions. \vec{x} defines the schema of \vec{c} , specifically, \vec{x} is a vector of names of each element of the tuple \vec{c} . It is not necessary that \vec{x} has the same schema as the given expression. Thus, we can evaluate $\text{dom}_{\vec{x}}$ for a broader range of \vec{x} and it is not restricted by the schema of the input query expression. In such cases the dom is infinite as the not presenting variables in the query can take any value.

For the comparison operator $(v_1 \theta v_2)$, where v_1 and v_2 are variables, we can compute the domain as follows:

$$\text{dom}_{\vec{x}}(v_1 \theta v_2) = \left\{ \vec{c} \mid \left(\left(\prod_i (x_i \leftarrow c_i) \right) (v_1 \theta v_2) \right) \neq 0 \right\} \quad (4)$$

The domain of a comparison is infinite. We can rewrite the above definition as a mathematical formula as follow:

$$\text{dom}_{\vec{x}}(v_1 \theta v_2) = \left\{ \vec{c} \mid \forall i, j : (v_1 = x_i \wedge v_2 = x_j) \Rightarrow c_i \theta c_j \right\} \quad (5)$$

For the join operator we can write:

$$\text{dom}_{\vec{x}}(q_1 \cdot q_2) = \{ \vec{c} \mid \vec{c} \in \text{dom}_{\vec{x}}(q_1) \wedge \vec{c} \in \text{dom}_{\vec{x}}(q_2) \} \quad (6)$$

while for the union operator the domain definition is very similar:

$$\text{dom}_{\vec{x}}(q_1 + q_2) = \{ \vec{c} \mid \vec{c} \in \text{dom}_{\vec{x}}(q_1) \vee \vec{c} \in \text{dom}_{\vec{x}}(q_2) \} \quad (7)$$

$$\text{dom}_{\vec{x}}(\text{constant}) = \{ \vec{c} \} \quad (8)$$

$$\text{dom}_{\vec{x}}(\text{variable}) = \{ \vec{c} \} \quad (9)$$

Finally, we can give a formalism for expressing the domain of a variable that will participate in an assignment operation:

$$\text{dom}_{\vec{x}}(v \leftarrow q_1) = \left\{ \vec{c} \mid \vec{c} \in \text{dom}_{\vec{x}}(q_1) \wedge \left(\forall i : (x_i = v) \Rightarrow \left((q_1 \cdot \prod_{j: x_j \neq v} (x_j = c_j)) = c_i \right) \right) \right\} \quad (10)$$

In fact using implication operator in the above definition allows us to extend the \vec{x} to whatever vector we want, as we already said the schema of \vec{x} is not necessarily the same as the schema of q . Here is a mathematical reformulation of (10):

$$\text{dom}_{\vec{x}}(v \leftarrow q_1) = \left\{ \vec{c} \mid \vec{c} \in \text{dom}_{\vec{x}}(q_1) \wedge (\forall i, j : (x_i = v = x_j) \Rightarrow (c_i = c_j)) \right\} \quad (11)$$

In Algorithm 1 we will take the domain of a domain. We use this for extending the domain regarding to a new schema vector \vec{x} . In other words we add some other variables to the vector and therefore compute the domain for each variable within that vector. Without lose of generality, we can consider a domain as a set with schema or a relation. So we can easily extend the domain of a set to a broader set of variables with the definition of the domain or a relation in (2). For the sake of completeness we define the domain of empty set as infinite:

$$\text{dom}_{\vec{x}}(\text{NULL}) = c \quad (12)$$

2 Computing the arities of the AGCA expressions

Having the definitions for the domains, we will try to give some insight regarding to the notion of arity. We will start with an example relation:

$$q = R(a, b) \cdot S(b, c)$$

the schema for q will have three variables (a, b, c) . The arity of a tuple is the number of its occurrences in the relation, in other words the order of multiplicity of that tuple. The arity of a tuple will increase or decrease if insertions or respectively deletions will be made to a relation. For example:

R	a	b	arity
	<1	<2>	1
	<1	<3>	2
	<3	<4>	1

Table 1: Relation R

S	b	c	arity
	<1	1>	1
	<2	2>	2
	<3	5>	2

Table 2: Relation S

$R \cdot S$	a	b	c	arity
	<1	2	2>	2
	<1	3	5>	4
	<3	4	* >	0
	< *	1	1>	0

Table 3: Relation $R \cdot S$

We need to maintain the maps for certain values as long as the arity of a tuple is greater then 0. If the arity of the tuple drops to 0 then the tuple will not be taken into consideration and therefore it can be eliminated from the domains of the maps.

We need to store the arities inside each map and also way to compute the arity of an AGCA expression. Since we substitute the subexpressions with maps, we can easily consider the relations as the maps without input variables. Also a map with input variables can be seen as a relation with a group-by clause. Input variable of a map bind some variables. Thus if we compute the map values by all different combinations of these variables, we will look up into these values and return the appropriate value according to the input variables.

We define the function *arity* which will be used to compute the arity of a tuple in a certain relation. The function will be defined on the relation and the tuple for which the multiplicity order is desired to be computed.

$$arity(\text{Relation } q, \text{Tuple } t) = \text{multiplicity order of } t \text{ in the } q = \pi_t(q)$$

where $\pi_t(q)$ means the projection of relation q for the tuple t . This function can be used for the computation of the arity of the expressions from the AGCA: $q ::= q \cdot q \mid q + q \mid q \theta t \mid t \leftarrow q \mid \text{constant} \mid \text{variable}$. However constants and variables can be eliminated from the computation because relations are of interest.

$$arity(q_1 \cdot q_2, t) = \sum_{\{t_1\} \bowtie \{t_2\} = t} arity(q_1, t_1) * arity(q_2, t_2) \quad (13)$$

$$arity(q_1 + q_2, t) = arity(q_1, t) + arity(q_2, t), \text{ where } \text{Schema}(q_1) \quad (14)$$

$$= \text{Schema}(q_2) \quad (15)$$

$$arity(v_1 \theta v_2, t) = \begin{cases} 1, & \text{if } v_1 \theta v_2 \text{ is true and } v_1 \in t \wedge v_2 \in t \\ 0, & \text{otherwise is false} \end{cases} \quad (16)$$

$$arity(v \leftarrow q, t) = \begin{cases} 0, & \text{if } t \setminus \{v\} \notin \text{dom}_{\text{Schema}(t)}(q) \\ 1, & \text{otherwise} \end{cases} \quad (17)$$

3 Domain computation

We are trying to compute the domain of variables that appear in queries. This computation is performed on the query decomposition tree. We can traverse the tree in post order, first visiting the leaves that are represented by some relations and afterwards visiting the parent nodes and combining the relations of the children nodes. Using this technique we are trying to compute the domains, but also the vector \vec{x} of all the variables defined in that query.

The intuition behind the algorithm is simple. For computing the domain of a node we first compute the domain of its left child and according to the operator of the node itself we decide how to send the information from left child to right child, and then the domain of the right child.

The algorithm needs as inputs the root of the tree and a structure that must be previously defined and returns a structure that contains the vector of variables(\vec{x}) and the domains for each variable defined in the vector(**dom**). The structure will look like:

```
struct {
x: the vector of all the variables
dom: the domains of all the variables
}
```

In Algorithm 1 we define function **computeDomain** which helps us to compute the domains for variables within a given query. When invoking the function, we pass as arguments of the function, the root of the decomposition tree

and a structure which will have the vector of variables and the domain of the variables as nil, $\text{computeDomain}(\text{root}, s)$.

Algorithm 1 $\text{computeDomains}(\text{node}, s)$

Input: node as the root of the tree to be traversed, and s as the structure

Output: a structure result that will contain the vector of variables \vec{x} and the domains of each variable $\text{dom}_{\vec{x}}(\text{query})$

```

1: if  $\text{node.type} = "+"$  then
2:    $s_1 \leftarrow \text{computeDomain}(\text{node.left}, s)$ 
3:    $s_2 \leftarrow \text{computeDomain}(\text{node.right}, s)$ 
4:    $\text{result}.\vec{x} \leftarrow s_1.\vec{x} \cup s_2.\vec{x}$  // this will compute the vector for all the
      variables, both from the right and left node
5:    $\text{result.dom} \leftarrow \text{dom}_{\text{result}.\vec{x}}(\text{node.left} + \text{node.right})$ 
6: else if  $\text{node.type} = "*" \text{ then}$ 
7:    $s_1 \leftarrow \text{computeDomain}(\text{node.left}, s)$ 
8:    $\text{result} \leftarrow \text{computeDomain}(\text{node.right}, s_1)$ 
9: else
10:   $\text{result}.\vec{x} \leftarrow s.\vec{x} \cup \{\text{all the variables of the node}\}$ 
11:   $\text{result.dom} \leftarrow \text{dom}_{\text{result}.\vec{x}}(s.\text{dom}) \cap \text{dom}_{\text{result}.\vec{x}}(\text{node})$ 
12: end if
13: if  $\text{node.map} = \text{true}$  then
14:    $\text{Maps}[\text{node}].\vec{x} \leftarrow \text{dom}_{\vec{x}}(s.\text{dom})$ 
15:    $\text{Maps}[\text{node}].\vec{y} \leftarrow \text{result}$ 
16: end if
17: return  $\text{result}$ 

```

The order of Algorithm 1 is $O(n \cdot P)$ where n is the number of nodes in the query decomposition tree and P is the time needed for any rule of dom 's computation, according to previous section. Clearly this algorithm visits each node only one time.

In Algorithm 1 we suppose that there are three types of nodes: union nodes that represent $q + q$ expressions, join nodes that represent $q * q$ expressions and map nodes, that will be represented by $M[\vec{x}][\vec{y}]$. In addition to nodes we have three types of leaves: simple relations $R(\vec{y})$, inequalities $v_1 \theta v_2$ and assignment relations $v \leftarrow q$. A union node or a join node will always has two children, therefore the line 9 will be executed when a leaf is visited. When a node that is represented by a map is encountered then the algorithm

should produce a set of variables and respectively the domains of each variable. Also in line 11 we use the definition of **dom** to extend its schema to a broader range according to what was said in the previous section, definition (2).

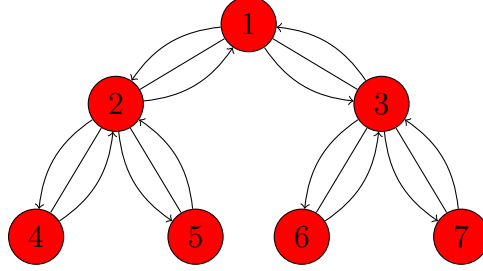


Figure 1: Tree traversal

The algorithm starts with node 1, which represents the root of the decomposition tree. It recursively goes through each of the leaves, therefore constructs the variable vector and also the domain for each variable.

It is known that the left most leaf will always be a relation that produces only output variables. Otherwise, if this condition is not sustained then it will contradict the rules imposed by the DBToaster calculus. Therefore node 4 will always be a simple relation that will give some output variables, which may help in future nodes if those nodes have input variables. If, for example node 5 will have input variables, then those variables will depend only on the output variables that node 4 is giving.

We have mentioned that besides map nodes, there are two more types of nodes: the union node and the join node. When a union node is met, then it is known that domains of variables will not pass over this operator and for that the variables and their domains, which were obtained in the parent, will be passed to each of this node's children, thus applying the distribution rule of join operator over the union operator. For example if we have $q1 * (q2 * q3 + q4 * q5) \equiv q1 * q2 * q3 + q1 * q4 * q5$, the variables and domain obtained for $q1$ will be passed both to $q2 * q3$ and $q3 * q4$, without modifying any of the domains, regarding that some domains may change on one branch, for example $q2 * q3$. For the other type of node, the join node, the passing of domains is allowed, and therefore a domain can be refined by the relation on the right side of the parent node.

Every time the algorithm encounters a map then it will create the domain for that specific map. Here we can make a difference between input and output variables, because the domain of input variables is dependent on the domains of output variables from the left side of the map's parent node, and the domain of the output variables are going to be produced by the children of the map's node. All this information will be stored in a global variable which can be easily accessed.

Theorem 1. *Algorithm 1 computes the domain of each variables and maps.*

Proof. The proof is by induction on the height of the tree. The theorem is held for a tree consisting of only one node. Since this node is a simple relation and the Algorithm 1 goes through lines 9-11. As we defined, the domain of any empty set is infinite(all possible values) so the domain is correctly computed in line 11.

Now suppose that the Algorithm 1 works for all the trees of a height less than k . We want to prove that it will work as well for trees of height k . Let T be a tree with height k . The root should have two children since the construction of the tree only creates nodes with leaves or two children [1]. The height of both of these children should be less than k . So we have correctly compute the domain of the left child and its input variable. Now we have two cases, if the root node is a union node or a join node. If this is a union node then we can compute the domain of the right child independently, since no information is passed over the “+” in DBToaster[1]. We compute the domain of the root according to definition (7). But if the root node is a join node, we have to pass the information from left child to the right child, since the right child may have some input variables which are defined in the left child. This is done in lines 6,7.

Thus in either of two cases we compute the domain of the root correctly and the theorem follows. \square

The algorithm presented above will compute the domains for each map that is met during the tree traversal. However, if, for example, we want to compute the domain of just one map, then the algorithm can easily be modified to accept this possibility. For example take the next figure 2.

Here we want to compute the domain for node 3. The algorithm will go exactly as the previous one, however the domain for the output variables can not be computed because the algorithm will not see any more children of the node represented by the map. Therefore, the use of this algorithm will

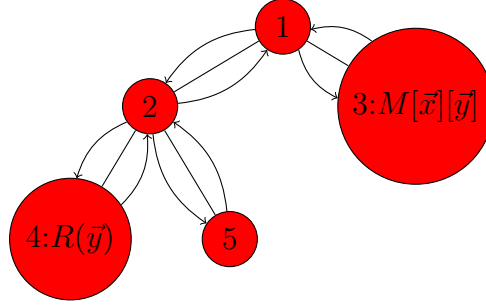


Figure 2: Partial tree seen by the algorithm

produce only the domains for the input variables, because the computation of the domains of the output variables will require the traversal of node's children.

For the algorithm to work, we must add a control variable which could be of boolean type. When invoking the function, the necessary arguments are going to be the root of the decomposition tree, the structure s , which will have the vector variable and the domains for each variable as nil , and the variable *stop* which will be transmitted by its address, not by the value, therefore any change to this variable will be permanently. At first the boolean variable will be *false*, but when the necessary map is encountered then the variable will become true, therefore stopping the process of computing the whole domain for each variable that might appear in the query expression.

Algorithm 2 computeDomains($node, s, stop$)

Input: $node$ as the root of the tree, s as the structure, $stop$ a boolean variable that will tell us when the specific map is met

Output: a structure $result$ that will contain the vector of variables \vec{x} and the domains of each variable $\text{dom}_{\vec{x}}(query)$

```
1: if stop=false then
2:   if  $node.type = "+"$  then
3:      $s_1 \leftarrow \text{computeDomain}(node.left, s, stop)$ 
4:      $s_2 \leftarrow \text{computeDomain}(node.right, s, stop)$ 
5:     if  $s_2! = nil$  then
6:        $result.\vec{x} \leftarrow s_1.\vec{x} \cup s_2.\vec{x}$  // this will compute the vector for all the
       variables, both from the right and left node
7:        $result.dom \leftarrow \text{dom}_{result.\vec{x}}(node.left + node.right)$ 
8:     else
9:        $result \leftarrow s_1$ 
10:    end if
11:  else if  $node.type = "*"$  then
12:     $s_1 \leftarrow \text{computeDomain}(node.left, s, stop)$ 
13:     $s_2 \leftarrow \text{computeDomain}(node.right, s_1, stop)$ 
14:    if  $s_2! = nil$  then
15:       $result \leftarrow s_2$ 
16:    else
17:       $result \leftarrow s_1$ 
18:    end if
19:  else if  $node.type = M[\vec{x}][\vec{y}]$  then
20:     $result.dom \leftarrow \text{dom}_{\vec{x}}(s.dom)$ 
21:     $result.\vec{x} \leftarrow \vec{x}$ 
22:     $stop = true$ 
23:  else
24:     $result.\vec{x} \leftarrow s.\vec{x} \cup \{\text{all the variables of the } node\}$ 
25:     $result.dom \leftarrow \text{dom}_{result.\vec{x}}(s.dom) \cap \text{dom}_{result.\vec{x}}(node)$ 
26:  end if
27: else
28:    $result = nil$ 
29: end if
30: return  $result$ 
```

4 Computing and maintaining the domains using the arities

The main goal is to maintain the domains of the variables of each map that will be present in query expression. We must offer a solution of maintaining these domains, see how the domains will increase or decrease if tuples are being added or deleted from the relations.

References

- [1] C. Koch, *Incremental Query Evaluation in a Ring of Databases*, preprint (2011).