# Aggregates in Logic

*Carlo Zaniolo,*

Computer Science Department

UCLA

# DBMSs' Support for Aggregates

- Only five aggregates in SQL2

- Vendors have added new builtins for data mining—e.g. rollups, datacubes, aggregates for time series

- these are never enough: **User-Defined Aggregates (UDAs)** are needed

- UDAs are in SQL3, but problems remain and not supported by vendors.

# More Flexible UDAs

---

- Currently aggregates are totally batch-oriented. For instance *Find cities where more than 20 employees live* can easily expressed with COUNT .

  But then SQL counts all the 10,000 employees in LA, before checking that this number is larger than 20

- A small percentage of data is often sufficient for a good estimate of averages. They propose the concept of on-line aggregation to solve this problem [Hellerstein, Hass and Wang]

- To solve this problem we introduce **early returns** in our UDAs.

- The Nonmonotonicity issue.

# NonMonotonicity

- COUNT: *find the suppliers who DO NOt supply red parts* $\leftrightarrow$ *find suppliers where the count of parts they supply is zero.*
  MAX and MIN: the highest paid employee is the one for which there is no employee with higher salary.

- Recursive queries are now supported in O/R Databases using techniques and semantics adapted from Deductive Databases: differential fixpoint, magic sets, stratified negation and aggregates.

- SQL queries **must be stratified w.r.t. negation and aggregates (SQL3)**.

# A Logical Recostruction of Aggregates

***

## Inductive Definition of Aggregates

1. $BASE$ For a singleton set: $count(\{x\}) = 1;\quad sum(\{x\}) = x;$ $max(\{x\}) = x$

2. $INDUCTION\ sum(S \sqcup \{x\}) = sum(S) + x;\quad count(S \sqcup \{x\}) = sum(S) + 1$
   $max(S \sqcup \{x\}) = if\, x > max(S)$ then $x$ else $max(S)$.

3. These computations can be easily expressed by logical rules (e.g., the *single* and *multi* rules in $\mathcal{LDL}++$)

4. Early returns and Final returns can also be expressed by rules: e.g., Avg= Sum/Count But we need to enumerate the elements of a set.

# Stable Models and Choice Models

- NonMonotonic Reasoning in AI, 30 years of progress: from circumscription to stable models.

- But nondetermism is also an essential facet of stable models. The following programs has dual models.

  a $\leftarrow \neg$b. b $\leftarrow \neg$a.

- Positive programs with choice can be rewritten as equivalent programs with negated goals displaying a multiplicity of stable models.

  advisor(St, Prof) $\leftarrow$ student(St, Maj),
  professor(Prof, Maj), choice((St), (Prof)).

- **Theorem** [PGZ]: Positive programs with choice define nondeterministic monotonic mappings.

# Aggregate Definition in $\mathcal{LDL}++$

Standard Average:

$$\mathtt{single}(\mathtt{avg}, \mathtt{Y}, (\mathtt{Y}, 1)).$$
$$\mathtt{multi}(\mathtt{avg}, \mathtt{Y}, (\mathtt{Sum}, \mathtt{Count}), (\mathtt{Sum} + \mathtt{Y}, \mathtt{Count} + 1)).$$
$$\mathtt{freturn}(\mathtt{avg}, \mathtt{Y}, (\mathtt{Sum}, \mathtt{Count}), \mathtt{Avg}) \qquad\qquad \leftarrow \mathtt{Avg} = \mathtt{Sum}/\mathtt{Count}.$$

On line average: returns a value every 100 samples.

$$\mathtt{ereturn}(\mathtt{avg}, \mathtt{X}, (\mathtt{Sum}, \mathtt{Count}), \mathtt{Avg}) \leftarrow \mathtt{Count} \bmod 100 = 0, \mathtt{Avg} = \mathtt{Sum}/\mathtt{Count}.$$

Using the aggregate remains the same:

$$\mathtt{p}(\mathtt{DeptNo}, \mathtt{avg}\langle\mathtt{Sal}\rangle) \leftarrow \mathtt{empl}(\mathtt{Ename}, \mathtt{Sal}, \mathtt{DeptNo}).$$

# Aggregates Formal Definition

$$p(\text{avg}\langle Y \rangle) \leftarrow d(Y).$$

We replace this by

$$p(Y) \leftarrow \text{results}(\text{avg}, Y).$$

where $\text{results}(\text{avg}, Y)$ is derived from $d(Y)$ by

- the chain rules,

- the **cagr** rules and

- the **return** rules.

# Aggregates: Definition

The chain rules are those with choice that place the elements of $d(Y)$ into an order-inducing chain.

Then, the `cagr` rules perform the inductive computation by calling the `single` and `multi` rules as follows:

$$\text{cagr}(\text{AgName}, Y, \text{New}) \leftarrow \quad \text{chain}(\text{nil}, Y), Y \neq \text{nil}, \text{single}(\text{avg}, Y, \text{New}).$$
$$\text{cagr}(\text{AgName}, Y2, \text{New}) \leftarrow \text{chain}(Y1, Y2), \text{cagr}(\text{AgName}, Y1, \text{Old}),$$
$$\text{multi}(\text{AgName}, Y2, \text{Old}, \text{New}).$$

Thus, the `cagr` rules are used to memorize the previous results, and to apply (i) `single` to the first element of $d(Y)$ (i.e., for the pattern `chain(nil, Y)`) and (ii) `multi` to the successive elements. The return

rules are as follows:

$$\text{results}(\text{AgName}, \text{Yield}) \leftarrow \text{chain}(\text{Y1}, \text{Y2}), \text{cagr}(\text{AgName}, \text{Y1}, \text{Old}),$$
$$\text{ereturn}(\text{AgName}, \text{Y2}, \text{Old}, \text{Yield}).$$
$$\text{results}(\text{AgName}, \text{Yield}) \leftarrow \text{chain}(\text{X}, \text{Y}), \neg\text{chain}(\text{Y}, \_),$$
$$\text{cagr}(\text{AgName}, \text{Y}, \text{Old}),$$
$$\text{freturn}(\text{AgName}, \text{Y}, \text{Old}, \text{Yield}).$$

Therefore we first compute `chain`, and then `cagr` that applies the `single` and `multi` to every element in the chain.

Concurrently, the first `results` rule produces all the results that can be generated by the application of `ereturn` to each element of the chain.

# Monotonic Aggregates

The final returns are computed by the second `results` rule that cals `freturn` once the last element in the chain (i.e., the element without successors) is detected.

This is the only rule using negation; in the absence of `freturn` this rule can be removed yielding a positive choice program that is monotonic!

Thus every aggregate with only early returns is *monotonic with respect to set-containment* and can be used freely in recursive rules. *monotonic aggregates*

To define a new aggregate, the user must write the `single`, `multi`, `ereturn` and `freturn` rules; the remaining rules are built in the system.

# Monotonic Aggregates: no `freturn`

---

Continuous count:

$$\text{single}(\text{mcount}, Y, 1).$$
$$\text{multi}(\text{mcount}, Y, \text{Old}, \text{New}) \leftarrow \quad \text{New} = \text{Old} + 1.$$
$$\text{ereturn}(\text{mcount}, Y, \text{Old}, \text{New}) \leftarrow \text{Old} = \text{nil}, \text{New} = 1$$
$$\text{ereturn}(\text{mcount}, Y, \text{Old}, \text{New}) \leftarrow \text{Old} \neq \text{nil}, \text{New} = \text{Old} + 1.$$

You can define msum in a similar fashion.

# mcount and msum

- Say that instead of `count` and `sum` we use `mcount`, `msum`, which returns a new partial count or sum for each new element in the set.

- Thus for a set of cardinality 5 `mcount` returns: $1, 2, 3, 4, 5$.

- If we add a new element to the set `mcount` returns: $1, 2, 3, 4, 5, 6$.

- `mcount` is monotonic and deterministic. But `msum` is a nondeterministic (i.e., multivalued) monotonic mapping

- New aggregates are conducive to more efficient algorithms.

# Monotonic Aggregates—Applications

The query,'*Find all departments with more than 7 employees*" can be expressed as follows:

$$\texttt{count\_emp}(\texttt{D\#}, \texttt{mcount}\langle\texttt{E\#}\rangle) \leftarrow \texttt{emp}(\texttt{E\#}, \texttt{Sal}, \texttt{D\#}).$$
$$\texttt{large\_dept}(\texttt{D\#}) \leftarrow \texttt{count\_emp}(\texttt{D\#}, \texttt{Count}), \texttt{Count} = 7.$$

*Find all departments with less than 7 employees*:

$$\texttt{small\_dept}(\texttt{D\#}, \texttt{Dname}) \leftarrow \texttt{dept}(\texttt{D\#}, \texttt{Dname}), \neg\texttt{large\_dept}(\texttt{D\#}).$$

# Monotonic Aggregates—Applications

Join the party: Some people will come to the party no matter what, and their names are stored in a **sure(Person)** relation. But others will join only after they know that at least $K = 3$ of their friends will be there. Here, $\texttt{friend}(\texttt{P}, \texttt{F})$ denotes that F is P's friend.

$$\texttt{willcome}(\texttt{P}) \leftarrow \qquad\qquad \texttt{sure}(\texttt{P}).$$
$$\texttt{willcome}(\texttt{P}) \leftarrow \qquad\qquad \texttt{c\_friends}(\texttt{P}, \texttt{K}), \texttt{K} >= 3.$$
$$\texttt{c\_friends}(\texttt{P}, \texttt{mcount}\langle \texttt{F} \rangle) \leftarrow \texttt{willcome}(\texttt{F}), \texttt{friend}(\texttt{P}, \texttt{F}).$$

sure(mark).
sure(tom).
sure(jane).

friend(jerry, mark).
friend(penny, mark).
friend(jerry, jane).
friend(penny, jane).
friend(jerry, penny).
friend(penny, tom).

Basic semi-naive computation yields:

willcome(mark).
willcome(tom).
willcome(jane).

c_friends(jerry, 1).
c_friends(penny, 1).
c_friends(jerry, 2).
c_friends(penny, 2).
c_friends(penny, 3).

willcome(penny).

c_friends(jerry, 3).

willcome(jerry).