SQL: Ordering Results, Duplicate Semantics and NULL Values Spring 2018

School of Computer Science University of Waterloo

Databases CS348

Ordering Results

- No particular ordering on the rows of a table can be assumed when queries are written. (This is important!)
- No particular ordering of rows of an intermediate result in the query can be assumed either.
- However, it is possible to order the final result of a query, using the ORDER BY clause at the end of the query.

General form:

ORDER BY
$$e_1[Dir_1], \ldots, e_k[Dir_k]$$

where Dir_i is either ASC or DESC; ASC is the default.

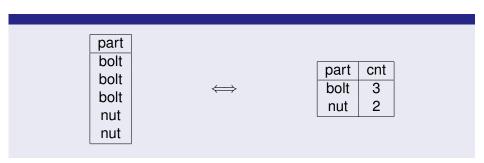
Example

List all authors in the database in ascending order of their name:

Again, the asc keyword is optional, and is assumed by default. A *descending* order is obtained with the desc keyword. Minor sorts, minor minor sorts, etc., can be added.

Multisets and Duplicates

- SQL has always had a MULTISET/BAG semantics rather than a SET semantics:
 - ⇒ SQL tables are **multisets** of tuples
 - ⇒ originally for efficiency reasons
- What does "allows duplicates" mean?



How does this impact Queries?

Example (Cheap Quantification-Projection)

EMP		
Name	Dept	
Bob	CS	
Sue	CS	
Fred	PMath	
Barb	Stats	
Jim	Stats	

 $\xrightarrow{\{y|\exists x.\mathsf{EMP}(x,y)\}}$

Dept
CS
CS
PMath
Stats
Stats

Dept cnt
CS 2
PMath 1
Stats 2

Range-restricted Queries for Multisets

Definition (Range restricted formulas with DISTINCT)

A formula (condition) φ is *range restricted* when, for φ_i that are also range restricted, φ has the form

$$\begin{array}{ll} R(x_{i_1},\ldots,x_{i_k}), \\ \varphi_1\wedge\varphi_2\,, \\ \varphi_1\wedge(x_i=x_j) & (\{x_i,x_j\}\cap FV(\varphi_1)\neq\emptyset), \\ \exists x_i.\varphi_1 & (x_i\in FV(\varphi_1)), \\ \varphi_1\vee\varphi_2 & (FV(\varphi_1)=FV(\varphi_2)), \\ \varphi_1\wedge\neg\varphi_2 & (FV(\varphi_2)=FV(\varphi_1)), \text{ or } \\ \mathsf{DISTINCT}(\varphi). \end{array}$$

Duplicates and Queries

How do we define what an answer to a query is now?

Ideas

- 1 A finite valuation can appear k times (k > 0) as a query answer.
- 2 The number of duplicates is a function of the numbers of duplicates in formulas.
- **3 DB**, θ , $k \models \varphi$ reads "finite valuation θ appears k times in φ 's answer".

Definition (Multiset Semantics for the Relational Calculus)

```
\begin{array}{lll} \mathbf{DB}, \theta, \mathbf{0} \models \varphi & \text{if} & \mathbf{DB}, \theta \not\models \varphi \\ \mathbf{DB}, \theta, \pmb{k} \models R(x_1, \dots, x_k) & \text{if} & (\theta(x_1), \dots, \theta(x_k)) \in \mathbf{R} & \pmb{k} \text{ times} \\ \mathbf{DB}, \theta, \pmb{m} \cdot \pmb{n} \models \varphi \land \psi & \text{if} & \mathbf{DB}, \theta, \pmb{m} \models \varphi \text{ and } \mathbf{DB}, \theta, \pmb{n} \models \psi \\ \mathbf{DB}, \theta, \pmb{m} \models \varphi \land (x_i = x_j) & \text{if} & \mathbf{DB}, \theta, \pmb{m} \models \varphi \text{ and } \theta(x_i) = \theta(x_j) \\ \mathbf{DB}, \theta, \sum_{v \in D} n_v \models \exists x. \varphi & \text{if} & \mathbf{DB}, \theta[x := v], n_v \models \varphi \\ \mathbf{DB}, \theta, \pmb{m} + \pmb{n} \models \varphi \lor \psi & \text{if} & \mathbf{DB}, \theta, \pmb{m} \models \varphi \text{ and } \mathbf{DB}, \theta, \pmb{n} \models \psi \\ \mathbf{DB}, \theta, \pmb{m} - \pmb{n} \models \varphi \land \neg \psi & \text{if} & \mathbf{DB}, \theta, \pmb{m} \models \varphi \text{ DB}, \theta, \pmb{n} \models \psi \text{ and } m > n \end{array}
```

Duplicates and SQL

Allowing duplicates leads to additional syntax.

- A duplicate elimination operator
 - ⇒ "SELECT DISTINCT x" v.s. "SELECT x" in SELECT-blocks
- MULTISET (BAG) operators
 - ⇒ equivalents of *set operations*
 - ⇒ but with multiset semantics.

Example

```
SQL> select r1.publication
  2 from wrote r1, wrote r2
  3 where r1.publication=r2.publication
  4
       and r1.author<>r2.author;
PUBLICAT
ChSa98
ChSa98
ChTo98
ChTo98
ChTo98a
ChTo98a
```

 \Rightarrow for publications with *n* authors we get $O(n^2)$ answers!

Bag Operations

- Bag union: UNION ALL
 - \Rightarrow additive union: bag containing all in Q_1 and Q_2 .
- Bag difference: EXCEPT ALL
 - ⇒ subtractive difference (monus):
 - \Rightarrow a bag all tuples in Q_1 for which there is no "matching" tuple in Q_2 .
- Bag intersection: INTERSECT ALL
 - \Rightarrow a bag of all tuples taking the maximal number common to Q_1 and Q_2

Example

```
SOL> ( select author
 2 from wrote, book
 3 where publication=pubid )
 4 union all
 5 ( select author
 6 from wrote, article
    where publication=pubid );
   AUTHOR
```

Summary

- SQL covered so far:
 - simple SELECT BLOCK
 - 2 set operations
 - 3 duplicates and multiset operations
 - 4 formulation of complex queries, nesting of queries, and views
 - 5 aggregation
- Note that duplicates in subqueries occurring in where clauses will not change the results computed by the top-level query, but that this is not true for subqueries in with or from clauses.

"Pure" SQL Equivalence, Revisited

Recall how nesting in the WHERE clause is syntactic sugar:

```
SELECT r.b

FROM r

WHERE r.a IN (
SELECT DISTINCT b
FROM s

FROM s
) AS s
)

WHERE r.a = s.b
```

Rewriting does not generally hold if DISTINCT is removed.

What is a "null" value?

Phone			
Name	Office	Home	
Joe	1234	3456	
Sue	1235	?	

- Sue doesn't have home phone (value inapplicable)
- Sue has home phone, but we don't know her number

(value unknown)

Value Inapplicable

- Essentially poor schema design.
- Better design:

Office Phone		
Name Office		
Joe	1234	
Sue 1235		

Home Phone		
Name Home		
Joe	Joe 3456	

- Queries should behave as if asked over the above decomposition.
 - ⇒ (relatively) easy to implement

Value Unknown

Idea

Unknown values can be replaced by any domain value (that satisfies integrity constraints).

⇒ many possibilities (possible worlds)

Phone		
Name	Office	Home
Joe	1234	3456
Sue	1235	?

Phone		
Name	Office	Home
Joe	1234	3456
Sue	1235	0000

:

Phone		
Name	Office	Home
Joe	1234	3456
Sue	1235	9999

Value Unknown and Queries

How do we answer queries?

Idea

Answers true in all possible worlds W of an incomplete D.

Certain Answer

$$Q(D) = \bigcap_{W \text{ world of } D} Q(W)$$

 \Rightarrow answer common to all possible worlds.

Is this (computationally) feasible?

⇒ NO (NP-hard to *undecidable* except in trivial cases)

SQL's solution: a (crude) approximation

What can we do with NULLs in SQL?

expressions

- general rule: a NULL as a parameter to an operation makes (should make) the result NULL
- $1 + \text{NULL} \rightarrow \text{NULL}$, 'foo' | | NULL $\rightarrow \text{NULL}$, etc.

predicates/comparisons

three-valued logic (crude approximation of "value unknown")

set operations

unique special value for duplicates

aggregate operations

doesn't "count" (i.e., "value inapplicable")

Comparisons Revisited

Idea

Comparisons with a NULL value return UNKNOWN

Example

$$egin{array}{lll} 1 &= 1 & & \mbox{TRUE} \\ 1 &= \mbox{NULL} & & \mbox{UNKNOWN} \\ 1 &= 2 & & \mbox{FALSE} \end{array}$$

Still short of proper logical behaviour:

$$x = 0 \lor x \neq 0$$

should be always true (no matter what x is, including NULL!), but...

UNKNOWN and Boolean Connectives

Idea

Boolean operations have to handle UNKNOWN

⇒ extended truth tables for Boolean connectives

\land	Т	U	F
Τ	Τ	U	F
U	U	U	F
F	F	F	F

... for tuples in which x is assigned the NULL value we get:

$$x=0 \lor x \neq 0
ightarrow$$
unknown \lor unknown $ightarrow$ unknown

which is not the same as TRUE.

UNKNOWN in WHERE Clauses

How is this used in a WHERE clause?

- Additional syntax IS TRUE, IS FALSE, and IS UNKNOWN

 ⇒ WHERE <cond> shorthand for WHERE <cond> IS TRUE
- Special comparison IS NULL

List all authors for which we don't know a URL of their home page:

Counting NULLS

How do NULLs interact with counting (and aggregates in general)?

■ count (URL) counts only non-NULL URL's

```
⇒ count (*) counts "rows"
```

Outer Join

Idea

Allow "NULL-padded" answers that "fail to satisfy" a conjunct in a conjunction

- Extension of syntax for the FROM clause
 - \Rightarrow FROM R <j-type> JOIN S ON C
 - \Rightarrow the <j-type> is one of FULL, LEFT, RIGHT, or INNER
- Semantics (for R(x, y), S(y, z), and C = (r.y = s.y)).
 - 1 $\{(x, y, z) : R(x, y) \land S(y, z)\}$
 - $(x, y, \text{NULL}) : R(x, y) \land \neg (\exists z. S(y, z))$ for LEFT and FULL
 - 3 $\{(\mathsf{NULL},y,z): S(y,z) \land \neg(\exists x.R(x,y))\}$ for RIGHT and FULL
 - \Rightarrow syntactic sugar for UNION ALL

Example

```
db2 => select aid, publication
db2 => from author left join wrote
db2 =>
                        on aid=author;
ATD
          PUBLICATION
          1 ChTo98
          1 ChTo98a
          1 Tom97
          2 ChTo98
          2 ChTo98a
          2 ChSa98
          3 ChSa98
  8 record(s) selected.
```

Counting with OJ

For every author count the number of publications:

```
db2 => select aid, count(publication) as pubs
db2 => from author left join wrote
dh2 =>
                       on aid=author
db2 => group by aid;
ATD
     PUBS
  4 record(s) selected.
```

Summary

- NULLs are necessary evil
 - ⇒ used to account for (small) irregularities in data
 - ⇒ should be used sparingly
- Can be always avoided
 - ⇒ however some of the solutions may be inefficient
- You can't escape NULLs in practice
 - ⇒ easy fix for blunders in schema design
 - $\Rightarrow \dots$ also due to schema evolution, etc.