

Relational Algebra

1

2



Relational Query Languages

- * Query languages: Allow manipulation and retrieval of data from a database.
- Relational model supports simple, powerful QLs:
 - Strong formal foundation based on logic.
 - Allows for much optimization.
- Query Languages != programming languages!
 - QLs not expected to be "Turing complete". e.g query language doesnt have looping
 - QLs not intended to be used for complex calculations.
 - QLs support easy, efficient access to large data sets.

Database Management Systems 3ed, R. Ramakrishnan and J. Gehrke

Formal Relational Query Languages

- Two mathematical Query Languages form the basis for "real" languages (e.g. SQL), and for implementation:
 - Relational Algebra: More operational, very useful for representing execution plans.
 - Relational Calculus: Lets users describe what they want, rather than how to compute it. (Nonoperational, <u>declarative</u>.)

Preliminaries

- * A query is applied to *relation instances*, and the result of a query is also a relation instance. QUERY COMPOSITION
 - Schemas of input relations for a query are fixed (but query will run regardless of instance!) input and output have the same schema
 - The schema for the result of a given query is also fixed! Determined by definition of query language constructs.
- Positional vs. named-field notation:
 - Positional notation easier for formal definitions, named-field notation more readable.
 - Both used in SQL

Database Management Systems 3ed, R. Ramakrishnan and J. Gehrke

1

Example Instances

R1	sid	<u>bid</u>	day	
	22	101	10/10/96	ĺ
	58	103	11/12/96	

- "Sailors" and "Reserves" relations for our examples. S1
- We'll use positional or named field notation, assume that names of fields in query results are 'inherited' from names of fields in query input relations.

sid	sname	rating	age
22	dustin	7	45.0
31	lubber	8	55.5
58	rusty	10	35.0

2	sid	sname	rating	age
	28	yuppy	9	35.0
	31	lubber	8	55.5
	44	guppy	5	35.0
	58	rusty	10	35.0

atabase Management Systems 3ed, R. Ramakrishnan and J. Gehr

5

Relational Algebra



- * Basic operations:
 - Selection (σ) Selects a subset of rows from relation.
 - <u>Projection</u> (π) Deletes unwanted columns from relation.
 - <u>Cross-product</u> (X) Allows us to combine two relations.
 - <u>Set-difference</u> (—) Tuples in reln. 1, but not in reln. 2.
 - <u>Union</u> (∪) Tuples in reln. 1 and in reln. 2.
- * Additional operations:
 - Intersection, <u>join</u>, division, renaming: Not essential, but (very!) useful.
- * Since each operation returns a relation, operations can be composed! (Algebra is "closed".)

 Database Management Systems 3ed, R. Ramadrishnan and I. Gehrke

relational algebra projection: use set that eliminates duplicates sname rating Projection yuppy sgl select: use bag that doesn't eliminate lubber 8 * Deletes attributes that are not in duplicates 5 guppy projection list. rustv Schema of result contains exactly $\pi_{sname,rating}(S2) =$ select distinct sname, rating from s2 the fields in the projection list, with the same names that they had in the (only) input relation. Projection operator has to age eliminate *duplicates*! (Why??) 35.0 • Note: real systems typically 55.5 don't do duplicate elimination unless the user explicitly asks dtp://Cate elimination is expensive (either via hash table or for it. (Why not?) because assert O(nlogn))

sid sname rating age Selection 28 35.0 yuppy 10 58 rusty 35.0 * Selects rows that satisfy $\sigma_{rating>8}$ (S2) selection condition. * No duplicates in result! (Why?) * Schema of result sname rating identical to schema of (only) input relation. yuppy * Result relation can be rusty 10 the input for another selection and projection cant switch the order relational algebra $\frac{\pi_{sname,rating}(\sigma_{rating})}{(doesn't have associativity)}$ operation! (Operator composition.) pi sname, rating (sigma age > 50 (s2)) 8 sigma age > 50(s2) (pi sname, rating)

Union, Intersection, Set-Difference 个情况, projection后的结果就没有age attribute,所以不可能select age > 50 sid sname rating 22 dustin 45.0 * All of these operations take two input relations, which 31 lubber 55.5 58 10 35.0 must be union-compatible: rusty 44 5 35.0 guppy Same number of fields. 28 35.0 yuppy Corresponding fields have the same type. $S1 \cup S2$ ❖ What is the *schema* of result? sid sname rating age sname rating age 31 lubber 8 55.5 dustin 45.0 35.0 rustv 10 $S1 \cap S2$ Dath demna en m tys em Se R. (an kriement) Cehrke

Cross-Product

- ❖ Each row of S1 is paired with each row of R1.
- * Result schema has one field per field of S1 and R1, with field names `inherited' if possible.
 - Conflict: Both S1 and R1 have a field called sid.

(sid)	sname	rating	age	(sid)	bid	day
22	dustin	7	45.0	22	101	10/10/96
22	dustin	7	45.0	58	103	11/12/96
31	lubber	8	55.5	22	101	10/10/96
31	lubber	8	55.5	58	103	11/12/96
58	rusty	10	35.0	22	101	10/10/96
58	rusty	10	35.0	58	103	11/12/96

Renaming operator: ρ (C(1→sid1,5→sid2), S1×R1)

all	combinations/permutations of tuples

10

rename column 1 as sid1 and column 5 as sid2 from the output of S1 X R1

Ioins

join operator

cross product operato

* Condition Join: $R \bowtie_{c} S = \sigma_{c} (R \times S)$

$$R \bowtie_{c} S = \sigma_{c}(R \times S)$$

(sid)	sname	rating	age	(sid)	bid	day
22	dustin	7	45.0	58	103	11/12/96
31	lubber	8	55.5	58	103	11/12/96

$$S1 \bowtie_{S1.sid < R1.sid} R1$$

- * Result schema same as that of cross-product.
- ❖ Fewer tuples than cross-product, might be able to compute more efficiently

condition join Sometimes called a theta-join if \$1 join \$1 with condition: self join

11

if s1 join r1 with equal operator condition: equi-join

natural join: equal join plus the attribute names are the same

join also follows by a projection while eliminate duplicate column (attribute)

Joins

* Equi-Join: A special case of condition join where the condition *c* contains only *equalities*.

sid	sname	rating	age	bid	day
22	dustin	7	45.0	101	10/10/96
58	rusty	10	35.0	103	11/12/96

$$S1 \bowtie_{sid} R1$$

- * Result schema similar to cross-product, but only one copy of fields for which equality is specified.
- * Natural Join: Equijoin on all common fields. When no condition is specified, tables will equijoin on all common fields

Division

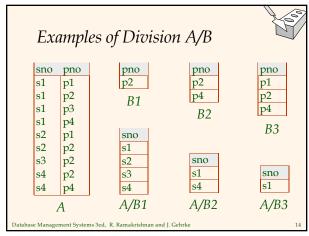


Not supported as a primitive operator, but useful for expressing queries like:

Find sailors who have reserved **all** boats.

- \star Let *A* have 2 fields, *x* and *y*; *B* have only field *y*:
 - $A/B = \{\langle x \rangle | \exists \langle x, y \rangle \in A \ \forall \langle y \rangle \in B \}$
 - i.e., *A/B* contains all *x* tuples (sailors) such that for *every y* tuple (boat) in *B*, there is an *xy* tuple in *A*.
 - *Or*: If the set of *y* values (boats) associated with an *x* value (sailor) in *A* contains all *y* values in *B*, the *x* value is in *A/B*.
- * In general, x and y can be any lists of fields; y is the list of fields in B, and $x \cup y$ is the list of fields of A.

13



14

Expressing A/B Using Basic Operators

- ❖ Division is not essential op; just a useful shorthand.
 - (Also true of joins, but joins are so common that systems implement joins specially.)
- ❖ Idea: For A/B, compute all x values that are not `disqualified' by some y value in B.
 - *x* value is *disqualified* if by attaching *y* value from *B*, we obtain an *xy* tuple that is not in *A*.

Disqualified x values: $\pi_{\chi}((\pi_{\chi}(A) \times B) - A)$

A/B: $\pi_{\chi}(A)$ – all disqualified tuples

Find names of sailors who've reserved boat #103

Solution 1: $\pi_{sname}((\sigma_{bid=103}^{} \text{Reserves}) \bowtie Sailors) \text{ Perform selects first (operation is cheaper)}$ solution 1 can be performed without writing temporary table to disk pipeline-able operators: ρ (Temp1, σ bid=103 Reserves) ρ (Temp2, Temp1 \bowtie Sailors) π_{sname} (Temp2)

Solution 3: $\pi_{sname}(\sigma_{bid=103}^{} \text{(Reserves)} \bowtie Sailors))$ Solution 3: $\pi_{sname}(\sigma_{bid=103}^{} \text{(Reserves)} \bowtie Sailors))$ Database Management Systems 3ed, R. Ramakrishnan and J. Gebrke

16

boats: bid, color, bname

Find names of sailors who've reserved a red boat

Information about boat color only available in Boats; so need an extra join:

\$\pi \text{sname}((\sigma_{color}='red', \begin{array}{cccc} Boats) \mathrew Reserves \mathrew Sailors) \\ \text{natural join} \end{array}\$

A more efficient solution:

\$\pi \text{sname}(\pi_{sid}(\pi_{bid}\sigma_{color}='red', \begin{array}{cccc} Boats) \mathrew Res) \mathrew Sailors) \\ \text{projection early on is cheaper because you are joining with less columns now many more tuples can fit into smaller number of disk pages

A query optimizer can find this, given the first solution!

17

Find sailors who've reserved a red or a green boat

* Can identify all red or green boats, then find sailors who've reserved one of these boats:

ρ (Tempboats, (σ color =' red' ∨ color =' green' Boats))

π sname (Tempboats ⋈ Reserves ⋈ Sailors)

natural join is associative (不管谁先谁后,结果一样) e.g 1+(2+3) = (1+2) + 3 join ordering (join tree ordering)

* Can also define Tempboats using union! (How?)

* What happens if ∨ is replaced by ∧ in this query?

Database Management Systems 3ed, R. Ramakrishnan and J. Gebrke

Find sailors who've reserved a red and a green boot

* Previous approach won't work! Must identify sailors who've reserved red boats, sailors who've reserved green boats, then find the intersection (note that sid is a key for Sailors):

 $\rho \ (Tempred, \pi_{sid}((\sigma_{color} = red, Boats)) \bowtie Reserves))$

 ρ (Tempgreen, π_{sid} ((σ_{color} -'green' Boats) \bowtie Reserves))

 $\pi_{sname}((Tempred \cap Tempgreen) \bowtie Sailors)$

join each type of boats with sailers first and then intersection ase Management Systems 3ed, R. Ramakrishnan and J. Gehrke

19 is red and green at the same time

Find the names of sailors who've reserved all boats

 Uses division; schemas of the input relations to / must be carefully chosen:

 $\rho \; (\textit{Tempsids}, (\pi_{\textit{sid,bid}}^{\mathsf{Reserves}}) \, / \, (\pi_{\textit{bid}}^{\mathsf{Boats}}))$

 π_{sname} (Tempsids \bowtie Sailors)

* To find sailors who've reserved all 'Interlake' boats:

..... $/\pi_{bid}(\sigma_{bname=Interlake})$ Boats)

20

Summary

- * The relational model has rigorously defined query languages that are simple and powerful.
- * Relational algebra is more operational; useful as internal representation for query evaluation plans.
- Several ways of expressing a given query; a query optimizer should choose the most efficient version.