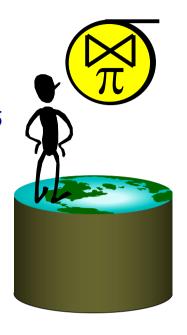
Relational Algebra

Fall 2016, Lecture 5

By relieving the brain of all unnecessary work, a good notation sets it free to concentrate on more advanced problems, and, in effect, *increases the mental power of the race*.

-- Alfred North Whitehead (1861 – 1947)



Formal Relational Query Languages

Two mathematical Query Languages form the basis for "real" languages (e.g. SQL), and for implementation:

<u>Relational Algebra</u>: More operational, very useful for representing execution plans.

Relational Calculus: Lets users describe what they want, rather than how to compute it. (Non-procedural, <u>declarative</u>.)

► Understanding Algebra (and Calculus) is key to understanding SQL, query processing!

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".
 - QLs not intended to be used for complex calculations.
 - QLs support easy, efficient access to large data sets.

Preliminaries

- A query is applied to *relation instances*, and the result of a query is also a relation instance.
 - Schemas of input relations for a query are fixed (but query will run over any legal instance)
 - The schema for the *result* of a given query is fixed.
 - It is determined by the definitions of the query language constructs.
- Positional vs. named-field notation:
 - Positional notation easier for formal definitions, named-field notation more readable.
 - Both used in SQL

Relational Algebra: 5 Basic Operations

- <u>Selection</u> (<u>o</u>) Selects a subset of **rows** from relation (horizontal).
- <u>Projection</u> (π) Retains only wanted **columns** from relation (vertical).
- <u>Cross-product</u> (x) Allows us to combine two relations.
- <u>Set-difference</u> (–) Tuples in r1, but not in r2.
- <u>Union</u> (\cup) Tuples in r1 and/or in r2.

Since each operation returns a relation, operations can be *composed!* (Algebra is "closed".)

Selection (σ) – Horizontal Restriction

- Selects rows that satisfy selection condition.
- Result is a relation.
 Schema of result is same as that of the input relation.

<u>sid</u>	sname	rating	age
28	yuppy	9	35.0
31	lubber	8	55.5
44	guppy	5	35.0
58	rusty	10	35.0

<u>sid</u>	sname	rating	age
28	yuppy	9	35.0
58	rusty	10	35.0

$$\sigma_{rating>8}$$
(S2)

Select all rows where the rating is larger than 8

Example Instances R1

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

<u>bid</u>	bname	color
101	Interlake	blue
102	Interlake	red
103	Clipper	green
104	Marine	red

Boats

<i>S</i> 1	sid	sname	rating	age
	22	dustin	7	45.0
	31	lubber	8	55.5
	58	rusty	10	35.0

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

Projection – Vertical Restriction

S2

- Examples: $\pi_{age}(S2)$; $\pi_{sname,rating}(S2)$
- Retains only attributes that are in the "projection list".
- Schema of result:
 - exactly the fields in the projection list, with the same names that they had in the input relation.
- Projection operator has to eliminate duplicates (How do they arise? Why remove them?)
 - Note: real systems typically don't do duplicate elimination unless the user explicitly asks for it. (Why not?)

(S2)

Projection

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

S2

sname	rating
yuppy	9
lubber	8
guppy	5
rusty	10

 $\pi_{sname,rating}(S2)$

age
35.0
55.5

$$\pi_{age}(S2)$$

Nesting Operators

- Result of a Relational Algebra Operator is a Relation, so...
- Can use as input to another Relational Algebra Operator

si	d	sname	rating	a	ge
2	8	yuppy	9	3.	5.0
3	_	lubber	-8	5	5.5
1	1		5	2	5.0
5	r R	guppy rusty	10	3	5.0
_		rasty	10		0.0

sname	rating
yuppy	9
rusty	10

 $\pi_{sname,rating}(\sigma_{rating>8}(S2))$

Union and Set-Difference

- All of these operations take two input relations, which must be *union-compatible*:
 - Same number of fields.
 - `Corresponding' fields have the same type.
- For which, if any, is duplicate elimination required?

Union

<u>sid</u>	sname	rating	age
22	dustin	7	45.0
31	lubber	8	55.5
58	rusty	10	35.0
S1			

sic	d sname	rating	age
22	dustin	7	45.0
31	lubber	8	55.5
58	rusty	10	35.0
44	guppy	5	35.0
28	yuppy	9	35.0

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

 $S1 \cup S2$

Set Difference

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

sid	sname	rating	age
22	dustin	7	45.0

$$S1-S2$$

S1

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

<u>sid</u>	sname	rating	age	
28	yuppy	9	35.0	
44	guppy	5	35.0	
S2-S1				

S2

Cross Product Example

S1	sid	sname	rating	age
	22	dustin	7	45.0
	31	lubber	8	55.5
	58	rusty	10	35.0

R1	sid	<u>bid</u>	day
	22	101	10/10/96
	58	103	11/12/96
			1

$$\rho(C(1 \rightarrow sid1, 5 \rightarrow sid2), S1 \times R1) =$$

sid1	sname	rating	age	sid2	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

Cross-Product

- S1 x R1: Each row of S1 paired with each row of R1.
 Q: How many rows in the result?
- Result schema has one field per field of S1 and R1, with field names `inherited' if possible.
 - May have a naming conflict: Both S1 and R1 have a field with the same name.
 - In this case, can use the *renaming operator*.

$$\rho$$
 (C(1 \rightarrow sid1,5 \rightarrow sid2), S1×R1)

Compound Operator: Intersection

- In addition to the 5 basic operators, there are several additional "Compound Operators"
 - These add no computational power to the language, but are useful shorthands.
 - Can be expressed solely with the basic ops.

Intersection takes two input relations, which must be *union-compatible*.

• Q: How to express it using basic operators?

$$R \cap S = R - (R - S)$$

Intersection

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

sid	sname	rating	age
31	lubber	8	55.5
58	rusty	10	35.0

S1

<u>sid</u>	sname	rating	age
28	yuppy	9	35.0
31	lubber	8	55.5
44	guppy	5	35.0
58	rusty	10	35.0

 $S1 \cap S2$

S2

Natural Join Example

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

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

S1

R1

R1⊳⊲S1 =

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

Compound Operator: Join (⋈)

- Joins are compound operators involving cross product, selection, and (sometimes) projection.
- Most common type of join is a "<u>natural join</u>" (often just called "join"). R ⋈S conceptually is:
 - Compute R X S
 - Select rows where attributes that appear in both relations have equal values
 - Project all unique attributes and one copy of each of the common ones.
- Note: Usually done much more efficiently than this.
- Useful for putting "normalized" relations back together.

Other Types of Joins

• Condition Join (or "theta-join"):

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

- Result schema same as that of cross-product.
- May have fewer tuples than cross-product.
- <u>Equi-Join</u>: Special case: condition c contains only conjunction of equalities.

"Theta" Join Example

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

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

S1

R1

$S1 \bowtie$	S1.sid < R1.sid $R1$	=
--------------	----------------------	---

(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

Division

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

Find sailors who have reserved all boats.

- Precondition: in A/B, the attributes in B must be included in the schema for A. Also, the result has attributes A-B.
 - SALES(supId, prodId);
 - PRODUCTS(prodId);
 - Relations SALES and PRODUCTS must be built using projections.
 - SALES/PRODUCTS: the ids of the suppliers supplying ALL products.

Division

- Assume
 - Relation R is defined over the attribute set A
 - Relation S is defined over the attribute set B
- Such that $B \subseteq A$ (B is a subset of A)
- Let C = A − B
- Division is defined as follows:
 - A relation over the attributes C that consists of the set of tuples from R that match the combination of every tuple in S.
- In other words, the result of R÷S consists of the restrictions of tuples in R to the attribute names unique to R, i.e., in the header of R but not in the header of S, for which it holds that all their combinations with tuples in S are present in R

Formally...

• A/B:

Let A have 2 fields, x and y; B have only field y:

A/B contains all x tuples such that for <u>every</u> y tuple in B, there is an xy tuple in A.]

$$A/B = \{\langle x \rangle | \forall \langle y \rangle \in B(\exists \langle x, y \rangle \in A) \}$$

- Why is this called division?
 - Answer: For all relations S and R it holds $S = (S \times R)/R$

More Examples of Division A/B

sno	pno	pno	pno	pn
s1	p1	p2	p2	p1
s1	p2	B1	\mathfrak{p}^{r}	p2
s1	р3	D1	B2	p4
s1	p4		DΔ	1.1
s2	p1	sno		B_{s}
s2	p2	s1		
s3	p2	s2	sno	
s4	p2	s3	s1	sn
s4	p4	s4	s4	s1
	4	A/B1	A/B2	\overline{A}

Note: For relation instances A and B, A/B is the largest relation instance Q such that $B \times Q \subseteq A$

Examples

Reserves	sid	<u>b</u>
	22	1

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

Sailors

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

Boats

bid	bname	color	
101	Interlake	Blue	
102	Interlake	Red	
103	Clipper	Green	
104	Marine	Red	

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:
$$\pi_{x}(A)$$
 – Disqualified x values

Find names of sailors who've reserved boat #103

• Solution 1:
$$\pi_{sname}((\sigma_{bid=103} \text{Reserves}) \bowtie Sailors)$$

• Solution 2:
$$\pi_{sname}(\sigma_{bid=103}(\text{Reserves}\bowtie Sailors))$$

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

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

Find names of sailors who've reserved a red boat

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

$$\pi_{sname}((\sigma_{color='red'}Boats) \bowtie Reserves \bowtie Sailors)$$

$$\pi_{\mathit{sname}}(\pi_{\mathit{sid}}((\pi_{\mathit{bid}}(\sigma_{\mathit{color}='\mathit{red}}, \mathit{Boats})) \bowtie \mathsf{Res}) \bowtie \mathit{Sailors})$$

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

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

	bname	color	
101	Interlake	blue	
102	Interlake	red	
103	Clipper	green	
104	Marine	red	

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

 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):

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

$$\rho \; (\textit{Tempgreen}, \pi_{\textit{sid}}((\sigma_{\textit{color} = '\textit{green'}}, \textit{Boats}) \bowtie \textit{Reserves}))$$

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

Find names of 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: $\rho \ (\textit{Tempboats}, (\sigma_{color} = \textit{red'} \lor color = \textit{green'} \ \textit{Boats}))$

 π_{sname} (Tempboats \bowtie Reserves \bowtie Sailors)

	Shame							
sid	sname	rating	age			<u>bid</u>	bname	colo
		_				101	Interlake	blue
22	dustin	7	45.0				Interlake	
31	lubber	Q	55.5					
-	lubbel	O					Clipper	greei
58	rusty	10	35.0			104	Marine	red
				1 • 1	1			

 sid
 bid
 day

 22
 101
 10/10/96

 58
 103
 11/12/96

Find the names of sailors who've reserved all boats

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

$$\rho \ (Tempsids, (\pi_{sid,bid}^{} Reserves) / (\pi_{bid}^{} Boats))$$
 $\pi_{sname} (Tempsids \bowtie Sailors)$

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

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