



Relational Algebra and Relational Calculus



Comp 521 - Files and Databases

Chapter 4



Fall 2014





Formal Query Languages

- What is the basis of Query Languages?
- Two formal Query Languages form the basis of "real" query languages (e.g. SQL):
 - <u>Relational Algebra</u>: Operational, it provides a recipe for evaluating the query. Useful for representing execution plans.
 - <u>Relational Calculus</u>: Lets users describe what they want, rather than how to compute it. (Non-operational, <u>declarative</u>.)







- Set of operands and operations that they are "closed" under all compositions
- Examples
 - Boolean algebra operands are the logical values True and False, and operations include AND(), OR(), NOT(), etc.
 - Integer algebra operands are the set of integers, operands include ADD(), SUB(), MUL(), NEG(), etc. many of which have special in-fix operator symbols (+,-,*,-)
- In our case operands are relations, what are the operators?





Example Instances

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

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

S1	<u>sid</u>	sname	rating	age
01	22	dustin	7	45.0
	31	lubber	8	55.5
	58	rusty	10	35.0

S2	sid	sname	rating	age
	28	yuppy	9	35.0
	31	lubber	8	55.5
	44	guppy	5	35.0
	58	rusty	10	35.0





Relational Algebra

Basic operations:

- <u>Selection</u> (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> (\cup) 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".)
 Comp 521 - Files and Databases
 Fall 2014





 Deletes attributes that are not in projection list.

Projection

- *Schema* of result contains exactly the fields in the projection list, with the same names that they had in the (only) input relation.
- Projection operator has to eliminate *duplicates*! (Why??)
 - Note: real systems typically don't do duplicate elimination unless the user explicitly asks for it. (Why not?)

<u>S1d</u>	sname	rating	age
28	yuppy	9	35.0
31	lubber	8	55.5
44	guppy	5	35.0
58	rusty	10	35.0
π		((S2)

sname,rating





Selection



- Selects rows that satisfy selection condition.
- No duplicates in result! (Why?)
- Schema of result
 identical to schema of
 (only) input relation.
- *Result* relation can be the *input* for another relational algebra operation! (*Operator composition*.)

sid	sname	rating	age
28	yuppy	9	35.0
31	lubber	8	55.5
44	guppy	5	35.0
58	rusty	10	35.0
	σ_{ratin}	$ng > 8^{(S2)}$	
	sname	rating	
	yuppy	9	
	rusty	10	
π		$(\sigma$	$(\mathbf{C2})$





Union, Intersection, Set-Difference

- All of these operations take two input relations, which must be <u>union-compatible</u>:
 - Same number of fields.
 - 'Corresponding' fields have the same type.
- What is the *schema* of result?

sid	sname	rating	age
22	dustin	7	45.0

sid	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

 $S1 \cup S2$

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

 $S1 \cap S2$





- 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:

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





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

(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

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

* *Result schema* same as that of cross-product.

- Fewer tuples than cross-product, might be able to compute more efficiently
- * Sometimes called a *theta-join*.



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
 (no labels on bowtie).





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

Find sailors who have reserved <u>all</u> boats.

- Let A have 2 fields, x and y; B have only field y:
 - $A/B = \{ \langle x \rangle \mid \exists \langle x, y \rangle \in A \forall \langle y \rangle \in B \}$
 - 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*.
 - 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*.





Examples of Division A/B





- Division is not essential; it's 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: $\pi_{\chi}(A)$ – disqualified tuples

Comp 521 – Files and Databases

Fall 2014





Relational Algebra Examples

- Assume the following extended schema:
 - Sailors(sid: integer, sname: string, rating: integer, age: real)
 - Reserves(sid: integer, bid: integer, day: date)
 - Boat(bid: integer, bname: string, bcolor: string)
- Objective: Write a relational algebra expression whose result instance satisfies the specified conditions
 - May not be unique
 - Some alternatives might be more efficient (in terms of time and/or space)





Names of sailors who've reserved boat #103

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

- * Solution 2: ρ (*Templ*, $\sigma_{bid=103}$ Reserves)
 - ρ (*Temp2*, *Temp1* \bowtie *Sailors*)

 π_{sname} (Temp2)

* Solution 3:
$$\pi_{sname}(\sigma_{bid=103}(\text{Reserves} \bowtie Sailors))$$





Names of sailors who've reserved a red boat

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

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

A more efficient solution:

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

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

Comp 521 – Files and Databases

Fall 2014





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:

 $\rho (Tempboats, (\sigma_{color ='red' \lor color ='green'} Boats))$

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

Can also define Tempboats using union! (How?)

* What happens if \vee is replaced by \wedge in this query?





Sailors who've reserved a red <u>and</u> 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):

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

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

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





Names of sailors who've reserved <u>all</u> boats

- Uses division; schemas of the input relations to / must be carefully chosen:
 - $\rho (Tempsids, (\pi_{sid, bid} \text{Reserves}) / (\pi_{bid} \text{Boats}))$ $\pi_{sname} (Tempsids \bowtie Sailors)$
- * To find sailors who've reserved all 'Interlake' boats:

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



Relational Calculus

- Comes in two flavors: <u>Tuple relational calculus</u> (TRC) and <u>Domain relational calculus</u> (DRC).
- Calculus has variables, constants, comparison ops, logical connectives and quantifiers.
 - <u>TRC</u>: Variables range over (i.e., get bound to) *tuples*.
 - <u>DRC</u>: Variables range over *domain elements* (= field values).
 - Both TRC and DRC are simple subsets of first-order logic.
- Expressions in the calculus are called *formulas with unbound formal variables*. An answer tuple is essentially an assignment of constants to these variables that make the formula evaluate to *true*.





- TRC and DRC are semantically similar
- In TRC, tuples share an equal status as variables, and field referencing can be used to select tuple parts
- In DRC, formal variables are explicit
- In the book you will find extensive discussions and examples of TRC Queries (Sections 4.3.1) and a lesser treatment of DRC.
- To even things out, in this lecture I will focus on DRC examples



Domain Relational Calculus

- * *Query* has the form: $\{<x1,x2,...,xn> | p(<x1,x2,...,xn>)\}$
- *Answer* includes all tuples <x1,x2,...,xn> that make the *formula* p(<x1,x2,...,xn>) *true*.
- * Formula is recursively defined, starting with simple atomic formulas (getting tuples from relations or making comparisons of values), and building bigger and better formulas using the logical connectives.





Atomic formula:

- $<x1,x2,...,xn> \in Rname$, or X op Y, or X op constant
- *op* is one of *<*,*>*,*=*,*≤*,*≥*,*≠*

Formula:

• an atomic formula, or



JX(p(X)) is read as "there exists a setting of the variable X such that p(X) is true". $\forall X(p(X))$ is read as "for all values of X, p(X)is true"

- $\neg p, p \land q, p \lor q$, where p and q are formulas, or
- **J**X(p(X)), where variable X is *free* in p(X), or
- $\forall X(p(X))$, where variable X is *free* in p(X)



- ✤ The use of quantifiers $\exists X$ and $\forall X$ in a formula is said to <u>bind</u> X.
 - A variable that is not bound is <u>free</u>.
- Let us revisit the definition of a query:

 $\{<x1,x2,...,xn> | p(<x1,x2,...,xn>)\}$

There is an important restriction: the variables x1, ..., xn that appear to the left of ' | ' must be the *only* free variables in the formula p(...).



Examples



Recall the example relations from last lecture

Sailors:

sid	sname	rating	age
22	Dustin	7	45.0
29	Brutus	1	33.0
31	Lubber	8	55.5
32	Andy	8	25.5
58	Rusty	10	35.0
64	Horatio	7	35.0
71	Zorba	10	16.0
74	Horatio	9	35.0
85	Art	3	25.5
95	Bob	3	63.5

Reservations:

sid	bid	day
22	101	10/10/98
22	102	10/10/98
22	103	10/8/98
22	104	10/7/98
31	102	11/10/98
31	103	11/6/98
31	104	11/12/98
64	101	9/5/98
64	102	9/8/98
74	103	9/8/98

Boats:

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



$\{ \langle I, N, T, A \rangle \mid \langle I, N, T, A \rangle \in Sailors \land T > 7 \}$

- ◆ The condition $\langle I, N, T, A \rangle \in Sailors$ binds the domain variables *I*, *N*, *T* and *A* to fields of any Sailors tuple.
- The term, <*I*,*N*,*T*,*A*>, to the left of '|' (which should be read as *such that*) says that every tuple, that satisfies *T* > 7 is in the answer.
- Modify this query to answer:
 - Find sailors who are older than 18 or have a rating under 9, and are called 'Joe'.



- ✤ Find all sailors with ratings above 7 $\{S \mid S \in Sailors \land S.rating > 7\}$
- ✤ Note, here S is a tuple variable

 $\{X \mid S \in Sailors \land S.rating > 7 \land X.name = S.name \land X.age = S.age \}$

Here X is a tuple variable with 2 fields (name, age). This query implicitly specifies projection (π) and renaming (ρ) relational algebra operators





Sailors rated > 7 who reserved boat #103

 $\{ \langle I, N, T, A \rangle \mid \langle I, N, T, A \rangle \in Sailors \land T \rangle 7 \land \\ \exists Ir, Br, D(\langle Ir, Br, D \rangle \in Reserves \land \\ Ir = I \land Br = 103) \}$

- ✤ We have used ∃ *Ir*, *Br*, D(...) as a shorthand for ∃ *Ir*(∃ *Br*(∃ D(...)))
- ♦ Note the use of ∃ to find a tuple in Reserves that 'joins with' (⋈) the Sailors tuples under consideration.



Find sailors rated > 7 who've reserved a red boat



 $\{ <I,N,T,A > | <I,N,T,A > \in Sailors \land T > 7 \land \\ \exists Ir, Br, D(<Ir, Br, D > \in Reserves \land Ir = I \land \\ \exists B, Bn, C(<B, Bn, C > \in Boats \land B = Br \land C = `red`)) \}$

- Observe how the parentheses control the scope of each quantifier's binding.
- This may look cumbersome, but with a good user interface, it is very intuitive. (MS Access, QBE)





$$\left\{ \left\langle N \right\rangle \middle| \exists I, T, A \left(\left\langle I, N, T, A \right\rangle \in Sailor \right) \\ \land \exists Ir, Br, D \left(\left\langle Ir, Br, D \right\rangle \in Reserves \land Ir = I \land Br = 103 \right) \right\}$$

- Note that only the *sname* field is retained in the answer and that only N is a free variable.
- A more compact version

$$\left\{ \left\langle N \right\rangle \middle| \exists I, T, A \left(\left\langle I, N, T, A \right\rangle \in Sailor \right) \right. \\ \wedge \exists D \left(\left\langle I, 103, D \right\rangle \in Reserves \right) \right\}$$

Comp 521 - Files and Databases

Fall 2014





Sailors who've reserved all boats

- Recall how queries of this type used of the "division" operator in relational algebra
- The trick is that we use "forall" quantification (∀) in place of "there exists" quantification (∃)
- Domains of variables are determined when they are bound
- Think of it as considering each variable's "domain" of independently in our substitution Comp 521 - Files and Databases

bid	bname	color
101	Interlake	blue
101	Interlake	red
101	Interlake	green
101	Clipper	blue
101	Clipper	red
101	Clipper	green
101	Marine	blue
101	Marine	red
101	Marine	green
102	Interlake	blue
	· ·	
104	Marine	green
104	marine	red



$$\left\{ \left\langle I, N, T, A \right\rangle \middle| \left\langle I, N, T, A \right\rangle \in Sailors \land \\ \forall B, BN, C \left(\neg \left(\left\langle B, BN, C \right\rangle \in Boats \right) \lor \\ \left(\exists Ir, Br, D \left(\left\langle Ir, Br, D \right\rangle \in Reserves \land I = Ir \land Br = B \right) \right) \right\}$$

Find all sailors I such that for each 3-tuple (B,BN,C) either it is not a tuple in Boats or there is a tuple in Reserves showing that sailor I has reserved it.







$$\begin{array}{l} \langle I, N, T, A \rangle | \langle I, N, T, A \rangle \in Sailors \land \\ \forall \langle B, BN, C \rangle \in Boats \\ (\exists \langle Ir, Br, D \rangle \in \operatorname{Reserves}(I = Ir \land Br = B)) \end{array}$$

Simpler notation, same query. (Much clearer!)
To find sailors who've reserved all red boats:

...
$$(C \neq 'red' \lor \exists \langle Ir, Br, D \rangle \in \operatorname{Reserves}(I = Ir \land Br = B))$$



It is possible to write syntactically correct calculus queries that have an infinite number of answers! Such queries are called <u>unsafe</u>.

• e.g.,
$$\left\{ < I, N, T, A > \left| < I, N, T, A > \notin Sailors \right\} \right\}$$

- It is known that every query that can be expressed in relational algebra can be expressed as a safe query in DRC / TRC; the converse is also true.
- Relational Completeness: Query language (e.g., SQL) can express every query that is expressible in relational algebra/calculus.





- Relational calculus is non-operational, and users define queries in terms of what they want, not in terms of how to compute it. (Declarativeness.)
- Algebra and safe calculus have same expressive power, leading to the notion of *relational completeness*.