

Lecture Notes: Relational Algebra

Det finns inget kapitel om relationsalgebra i kursen. Jag hade först tänkt ha med ett, men relationsalgebra passar inte riktigt i en grundkurs som den här. I stället finns en kort förklaring [i ordlistan](#), och för den som vill läsa mer finns dessutom dessa föreläsningssanteckningar på engelska.

What? Why?

- Similar to normal algebra (as in $2+3*x-y$), except we use relations as values instead of numbers, and the operations and operators are different.
- Not used as a query language in actual DBMSs. (SQL instead.)
- The inner, lower-level operations of a relational DBMS are, or are similar to, relational algebra operations. We need to know about relational algebra to understand query execution and optimization in a relational DBMS.
- Some advanced SQL queries requires explicit relational algebra operations, most commonly *outer join*.
- Relations are seen as *sets of tuples*, which means that no duplicates are allowed. SQL behaves differently in some cases. Remember the SQL keyword *distinct*.
- SQL is *declarative*, which means that you tell the DBMS *what* you want, but not *how* it is to be calculated. A C++ or Java program is *procedural*, which means that you have to state, step by step, exactly how the result should be calculated. Relational algebra is (more) procedural than SQL. (Actually, relational algebra is mathematical expressions.)

Det här är ett avsnitt i en [webbkurs om databaser](#) som finns fritt tillgänglig på adressen <http://www.databasteknik.se/webbkursen/>. Senaste ändring: 26 juni 2008.

Av [Thomas Padron-McCarthy](#). Copyright, alla rättigheter reserverade, osv. Skicka gärna kommentarer till webbkursen@databasteknik.se.

Set operations

Relations in relational algebra are seen as sets of tuples, so we can use basic set operations.

Review of concepts and operations from set theory

- set
- element
- no duplicate elements (but: multiset = bag)
- no order among the elements (but: ordered set)
- subset
- proper subset (with fewer elements)
- superset
- union
- intersection
- set difference
- cartesian product

Projection

Example: The table **E** (for **EMPLOYEE**)

nr	name	salary
1	John	100
5	Sarah	300
7	Tom	100

SQL	Result	Relational algebra								
<pre>select salary from E</pre>	<table><tr><th>salary</th></tr><tr><td>100</td></tr><tr><td>300</td></tr></table>	salary	100	300	$\text{PROJECT}_{\text{salary}}(E)$					
salary										
100										
300										
<pre>select nr, salary from E</pre>	<table><tr><th>nr</th><th>salary</th></tr><tr><td>1</td><td>100</td></tr><tr><td>5</td><td>300</td></tr><tr><td>7</td><td>100</td></tr></table>	nr	salary	1	100	5	300	7	100	$\text{PROJECT}_{\text{nr, salary}}(E)$
nr	salary									
1	100									
5	300									
7	100									

Note that there are no duplicate rows in the result.

Selection

The same table **E** (for **EMPLOYEE**) as above.

SQL	Result	Relational algebra									
select * from E where salary < 200	<table> <tr><td>nr</td><td>name</td><td>salary</td></tr> <tr><td>1</td><td>John</td><td>100</td></tr> <tr><td>7</td><td>Tom</td><td>100</td></tr> </table>	nr	name	salary	1	John	100	7	Tom	100	$\text{SELECT}_{\text{salary} < 200}(E)$
nr	name	salary									
1	John	100									
7	Tom	100									
select * from E where salary < 200 and nr >= 7	<table> <tr><td>nr</td><td>name</td><td>salary</td></tr> <tr><td>7</td><td>Tom</td><td>100</td></tr> </table>	nr	name	salary	7	Tom	100	$\text{SELECT}_{\text{salary} < 200 \text{ and nr} \geq 7}(E)$			
nr	name	salary									
7	Tom	100									

Note that the select operation in relational algebra has nothing to do with the SQL keyword select. Selection in relational algebra returns those tuples in a relation that fulfil a condition, while the SQL keyword select means "here comes an SQL statement".

Relational algebra expressions

SQL	Result	Relational algebra						
<pre>select name, salary from E where salary < 200</pre>	<table><tr><th>name</th><th>salary</th></tr><tr><td>John</td><td>100</td></tr><tr><td>Tom</td><td>100</td></tr></table>	name	salary	John	100	Tom	100	<p>PROJECT_{name, salary} (SELECT_{salary < 200}(E))</p> <p><i>or, step by step, using an intermediate result</i></p> <p>Temp <- SELECT_{salary < 200}(E) Result <- PROJECT_{name, salary}(Temp)</p>
name	salary							
John	100							
Tom	100							

Notation

The operations have their own symbols. The symbols are hard to write in HTML that works with all browsers, so I'm writing **PROJECT** etc here. The real symbols:

Operation	My HTML	Symbol
Projection	PROJECT	π

Operation	My HTML	Symbol
Cartesian product	X	\times

Selection	SELECT	σ
Renaming	RENAME	ρ
Union	UNION	\cup
Intersection	INTERSECTION	\cap
Assignment	<-	\leftarrow

Join	JOIN	\bowtie
Left outer join	LEFT OUTER JOIN	\bowtie_{L}
Right outer join	RIGHT OUTER JOIN	\bowtie_{R}
Full outer join	FULL OUTER JOIN	\bowtie_{F}
Semijoin	SEMIJOIN	\ltimes

Example: The relational algebra expression which I would here write as

PROJECT_{Namn} (**SELECT**_{Medlemsnummer < 3} (**Medlem**))

should actually be written

$\pi_{\text{Namn}} (\sigma_{\text{Medlemsnummer} < 3} (\text{Medlem}))$

Cartesian product

The *cartesian product* of two tables combines each row in one table with each row in the other table.

Example: The table **E** (for **EMPLOYEE**)

enr	ename	dept
1	Bill	A
2	Sarah	C
3	John	A

Example: The table **D** (for **DEPARTMENT**)

dnr	dname
A	Marketing
B	Sales
C	Legal

SQL	Result					Relational algebra
select * from E, D	enr	ename	dept	dnr	dname	E X D
	1	Bill	A	A	Marketing	
	1	Bill	A	B	Sales	
	1	Bill	A	C	Legal	
	2	Sarah	C	A	Marketing	
	2	Sarah	C	B	Sales	
	2	Sarah	C	C	Legal	
	3	John	A	A	Marketing	
	3	John	A	B	Sales	

3	John	A	C	Legal
---	------	---	---	-------

- Seldom useful in practice.
- Usually an error.
- Can give a huge result.

Join (sometimes called "inner join")

The cartesian product example above combined each employee with each department. If we only keep those lines where the **dept** attribute for the employee is equal to the **dnr** (the department number) of the department, we get a nice list of the employees, and the department that each employee works for:

SQL	Result	Relational algebra																				
<pre>select * from E, D where dept = dnr</pre>	<table><tr><th>enr</th><th>ename</th><th>dept</th><th>dnr</th><th>dname</th></tr><tr><td>1</td><td>Bill</td><td>A</td><td>A</td><td>Marketing</td></tr><tr><td>2</td><td>Sarah</td><td>C</td><td>C</td><td>Legal</td></tr><tr><td>3</td><td>John</td><td>A</td><td>A</td><td>Marketing</td></tr></table>	enr	ename	dept	dnr	dname	1	Bill	A	A	Marketing	2	Sarah	C	C	Legal	3	John	A	A	Marketing	<p>SELECT_{dept = dnr} (E X D)</p> <p><i>or, using the equivalent join operation</i></p> <p>E JOIN_{dept = dnr} D</p>
enr	ename	dept	dnr	dname																		
1	Bill	A	A	Marketing																		
2	Sarah	C	C	Legal																		
3	John	A	A	Marketing																		

- A very common and useful operation.
- Equivalent to a cartesian product followed by a select.
- Inside a relational DBMS, it is usually much more efficient to calculate a join directly, instead of calculating a cartesian product and then throwing away most of the lines.
- Note that the same SQL query can be translated to several different relational algebra expressions, which all give the same result.
- If we assume that these relational algebra expressions are executed, inside a relational DBMS which uses relational algebra operations as its lower-level internal operations, different relational algebra expressions can take very different time (and memory) to execute.

Natural join

A normal inner join, but using the join condition that columns with the same names should be equal. Duplicate columns are removed.

Renaming tables and columns

Example: The table **E** (for **EMPLOYEE**)

nr	name	dept
1	Bill	A
2	Sarah	C
3	John	A

Example: The table **D** (for **DEPARTMENT**)

nr	name
A	Marketing
B	Sales
C	Legal

We want to join these tables, but:

- Several columns in the result will have the same name (**nr** and **name**).
- How do we express the join condition, when there are two columns called **nr**?

Solutions:

- Rename the attributes, using the *rename* operator.
- Keep the names, and prefix them with the table name, as is done in SQL. (This is somewhat unorthodox.)

SQL	Result	Relational algebra																				
<pre>select * from E as E(enr, ename, dept), D as D(dnr, dname) where dept = dnr</pre>	<table><tr><th>enr</th><th>ename</th><th>dept</th><th>dnr</th><th>dname</th></tr><tr><td>1</td><td>Bill</td><td>A</td><td>A</td><td>Marketing</td></tr><tr><td>2</td><td>Sarah</td><td>C</td><td>C</td><td>Legal</td></tr><tr><td>3</td><td>John</td><td>A</td><td>A</td><td>Marketing</td></tr></table>	enr	ename	dept	dnr	dname	1	Bill	A	A	Marketing	2	Sarah	C	C	Legal	3	John	A	A	Marketing	$(\text{RENAME}_{(\text{enr}, \text{ename}, \text{dept})}(\text{E}))$ $\text{JOIN}_{\text{dept} = \text{dnr}} (\text{RENAME}_{(\text{dnr}, \text{dname})}(\text{D}))$
enr	ename	dept	dnr	dname																		
1	Bill	A	A	Marketing																		
2	Sarah	C	C	Legal																		
3	John	A	A	Marketing																		
<pre>select * from E, D where dept = D.nr</pre>	<table><tr><th>nr</th><th>name</th><th>dept</th><th>nr</th><th>name</th></tr><tr><td>1</td><td>Bill</td><td>A</td><td>A</td><td>Marketing</td></tr><tr><td>2</td><td>Sarah</td><td>C</td><td>C</td><td>Legal</td></tr><tr><td>3</td><td>John</td><td>A</td><td>A</td><td>Marketing</td></tr></table>	nr	name	dept	nr	name	1	Bill	A	A	Marketing	2	Sarah	C	C	Legal	3	John	A	A	Marketing	$\text{E JOIN}_{\text{dept} = \text{D.nr}} \text{D}$
nr	name	dept	nr	name																		
1	Bill	A	A	Marketing																		
2	Sarah	C	C	Legal																		
3	John	A	A	Marketing																		

You can use another variant of the renaming operator to change the name of a table, for example to change the name of **E** to **R**. This is necessary when joining a table with itself (see below).

$\text{RENAME}_R(\text{E})$

A third variant lets you rename both the table and the columns:

$\text{RENAME}_{R(\text{enr}, \text{ename}, \text{dept})}(\text{E})$

Aggregate functions

Example: The table **E** (for **EMPLOYEE**)

nr	name	salary	dept
1	John	100	A
5	Sarah	300	C
7	Tom	100	A
12	Anne	null	C

SQL	Result	Relational algebra		
<pre>select sum(salary) from E</pre>	<table><tr><th>sum</th></tr><tr><td>500</td></tr></table>	sum	500	$F_{\text{sum(salary)}}(E)$
sum				
500				

Note:

- Duplicates are not eliminated.
- **Null** values are ignored.

SQL	Result	Relational algebra
	Result:	

select count(salary) from E	count 3	$F_{\text{count(salary)}}(E)$
select count(distinct salary) from E	Result: count 2	$F_{\text{count(salary)}}(\mathbf{PROJECT}_{\text{salary}}(E))$

You can calculate aggregates "grouped by" something:

SQL	Result	Relational algebra
select sum(salary) from E group by dept	dept sum A 200 C 300	$\text{dept } F_{\text{sum(salary)}}(E)$

Several aggregates simultaneously:

SQL	Result	Relational algebra
select sum(salary), count(*) from E group by dept	dept sum count A 200 2 C 300 1	$\text{dept } F_{\text{sum(salary), count(*)}}(E)$

Standard aggregate functions: sum, count, avg, min, max

Hierarchies

Example: The table **E** (for **EMPLOYEE**)

nr	name	mgr
1	Gretchen	null
2	Bob	1
5	Anne	2
6	John	2
3	Hulda	1
4	Hjalmar	1
7	Usama	4

Going up in the hierarchy *one level*: What's the name of John's boss?

SQL	Result	Relational algebra
select b.name from E p, E b where p.mgr = b.nr and p.name = "John"	name Bob	$\mathbf{PROJECT}_{\text{b.name}} ([\text{SELECT}_{\text{pname} = \text{"John"}} (\mathbf{RENAME}_{\text{P}(\text{pnr}, \text{pname}, \text{pmgr})}(E))] \text{ JOIN}_{\text{pmgr} = \text{b.nr}} [\mathbf{RENAME}_{\text{B}(\text{bnr}, \text{bname}, \text{bmgr})}(E)])$ <i>or, in a less wide-spread notation</i> $\mathbf{PROJECT}_{\text{b.name}} ([\text{SELECT}_{\text{name} = \text{"John"}} (\mathbf{RENAME}_{\text{P}}(E))] \text{ JOIN}_{\text{p.mgr} = \text{b.nr}} [\mathbf{RENAME}_{\text{B}}(E)])$ <i>or, step by step</i>

```

P <- RENAMEP(pnr, pname, pmgr)(E)
B <- RENAMEB(bnr, bname, bmgr)(E)
J <- SELECTname = "John"(P)
C <- J JOINpmgr = bnr B
R <- PROJECTbname(C)

```

Notes about renaming:

- We are joining **E** with itself, both in the SQL query and in the relational algebra expression: it's like joining two tables with the same name and the same attribute names.
- Therefore, some renaming is required.
- **RENAME**_P(E) **JOIN**... **RENAME**_B(E) is a start, but then we still have the same attribute names.

Going up in the hierarchy *two levels*: What's the name of John's boss' boss?

SQL	Result	Relational algebra		
<pre>select ob.name from E p, E b, E ob where b.mgr = ob.nr where p.mgr = b.nr and p.name = "John"</pre>	<table><tr><td>name</td></tr><tr><td>Gretchen</td></tr></table>	name	Gretchen	<p>PROJECT_{ob.name} (((SELECT_{name = "John"}(RENAME_P(E))) JOIN_{p.mgr = b.nr} [RENAME_B(E)]) JOIN_{b.mgr = ob.nr} [RENAME_{OB}(E)])</p> <p><i>or, step by step</i></p> <p>P <- RENAME_{P(pnr, pname, pmgr)}(E) B <- RENAME_{B(bnr, bname, bmgr)}(E) OB <- RENAME_{OB(obnr, obname, obmgr)}(E) J <- SELECT_{name = "John"}(P) C1 <- J JOIN_{pmgr = bnr} B C2 <- C1 JOIN_{bmgr = bbnr} OB R <- PROJECT_{obname}(C2)</p>
name				
Gretchen				

Recursive closure

Both one and two levels up: What's the name of John's boss, *and* of John's boss' boss?

SQL	Result	Relational algebra			
(select b.name ...) union (select ob.name ...)	<table><tr><th>name</th></tr><tr><td>Bob</td></tr><tr><td>Gretchen</td></tr></table>	name	Bob	Gretchen	(...) UNION (...)
name					
Bob					
Gretchen					

Recursively: What's the name of *all* John's bosses? (One, two, three, four or more levels.)

- Not possible in (conventional) relational algebra, but a special operation called **transitive closure** has been proposed.
- Not possible in (standard) SQL (SQL2), but in SQL3, and using SQL + a host language with loops or recursion.

Outer join

Example: The table **E** (for **EMPLOYEE**)

enr	ename	dept
1	Bill	A
2	Sarah	C
3	John	A

Example: The table **D** (for **DEPARTMENT**)

dnr	dname
A	Marketing
B	Sales
C	Legal

List each employee together with the department he or she works at:

SQL	Result	Relational algebra																				
<pre>select * from E, D where edept = dnr</pre> <p><i>or, using an explicit join</i></p> <pre>select * from (E join D on edept = dnr)</pre>	<table><tr><th>enr</th><th>ename</th><th>dept</th><th>dnr</th><th>dname</th></tr><tr><td>1</td><td>Bill</td><td>A</td><td>A</td><td>Marketing</td></tr><tr><td>2</td><td>Sarah</td><td>C</td><td>C</td><td>Legal</td></tr><tr><td>3</td><td>John</td><td>A</td><td>A</td><td>Marketing</td></tr></table>	enr	ename	dept	dnr	dname	1	Bill	A	A	Marketing	2	Sarah	C	C	Legal	3	John	A	A	Marketing	$E \text{ JOIN}_{\text{edept} = \text{dnr}} D$
enr	ename	dept	dnr	dname																		
1	Bill	A	A	Marketing																		
2	Sarah	C	C	Legal																		
3	John	A	A	Marketing																		

No employee works at department B, Sales, so it is not present in the result. This is probably not a problem in this case. But what if we want to know the number of employees at each department?

SQL	Result	Relational algebra									
<pre>select dnr, dname, count(*) from E, D where edept = dnr group by dnr, dname or, using an explicit join select dnr, dname, count(*) from (E join D on edept = dnr) group by dnr, dname</pre>	<table> <tr> <th>dnr</th><th>dname</th><th>count</th></tr> <tr> <td>A</td><td>Marketing</td><td>2</td></tr> <tr> <td>C</td><td>Legal</td><td>1</td></tr> </table>	dnr	dname	count	A	Marketing	2	C	Legal	1	$\text{dnr, dname } F_{\text{count}(*)}(E \text{ JOIN}_{\text{edept} = \text{dnr}} D)$
dnr	dname	count									
A	Marketing	2									
C	Legal	1									

No employee works at department B, Sales, so it is not present in the result. It disappeared already in the join, so the aggregate function never sees it. But what if we want it in the result, with the right number of employees (zero)?

Use a *right outer join*, which keeps all the rows from the right table. If a row can't be connected to any of the rows from the left table according to the join condition, **null** values are used:

SQL	Result					Relational algebra
select * from (E right outer join D on edept = dnr)	enr	ename	dept	dnr	dname	E RIGHT OUTER JOIN _{edept = dnr} D
	1	Bill	A	A	Marketing	
	2	Sarah	C	C	Legal	
	3	John	A	A	Marketing	

	<table><tr><td>null</td><td>null</td><td>null</td><td>B</td><td>Sales</td></tr></table>	null	null	null	B	Sales								
null	null	null	B	Sales										
<pre>select dnr, dname, count(*) from (E right outer join D on edept = dnr) group by dnr, dname</pre>	<table><tr><td>dnr</td><td>dname</td><td>count</td></tr><tr><td>A</td><td>Marketing</td><td>2</td></tr><tr><td>B</td><td>Sales</td><td>1</td></tr><tr><td>C</td><td>Legal</td><td>1</td></tr></table>	dnr	dname	count	A	Marketing	2	B	Sales	1	C	Legal	1	dnr, dname ^F _{count(*)} (E RIGHT OUTER JOIN _{eddept = dnr} D)
dnr	dname	count												
A	Marketing	2												
B	Sales	1												
C	Legal	1												
<pre>select dnr, dname, count(enr) from (E right outer join D on edept = dnr) group by dnr, dname</pre>	<table><tr><td>dnr</td><td>dname</td><td>count</td></tr><tr><td>A</td><td>Marketing</td><td>2</td></tr><tr><td>B</td><td>Sales</td><td>0</td></tr><tr><td>C</td><td>Legal</td><td>1</td></tr></table>	dnr	dname	count	A	Marketing	2	B	Sales	0	C	Legal	1	dnr, dname ^F _{count(enr)} (E RIGHT OUTER JOIN _{eddept = dnr} D)
dnr	dname	count												
A	Marketing	2												
B	Sales	0												
C	Legal	1												

Join types:

- **JOIN** = "normal" join = inner join
- **LEFT OUTER JOIN** = left outer join
- **RIGHT OUTER JOIN** = right outer join
- **FULL OUTER JOIN** = full outer join

Outer union

Outer union can be used to calculate the union of two relations that are *partially union compatible*. Not very common.

Example: The table **R**

A	B
1	2
3	4

Example: The table **S**

B	C
4	5
6	7

The result of an outer union between R and S:

A	B	C
1	2	null
3	4	5
null	6	7

Division

Who works on (at least) *all* the projects that Bob works on?

Semijoin

A join where the result only contains the columns from one of the joined tables. Useful in distributed

databases, so we don't have to send as much data over the network.

Update

To update a named relation, just give the variable a new value. To add all the rows in relation **N** to the relation **R**:

R <- R UNION N

[Webbkursen om databaser](#) av [Thomas Padron-McCarthy](#).