

The Relational Algebra

Dirk Van Gucht¹

¹Indiana University

February 5, 2019

Outline

- Motivation
- Syntax of Relational Algebra (RA)
- Semantics of RA in SQL
- Expressing queries in RA

Relational Algebra (Motivation)

- SQL is a declarative language to specify queries over relational databases
- RA is a procedural language wherein algebraic expressions specify queries
- RA expressions provide procedures (algorithms) to evaluate queries
 - Examining these procedure provides insights about the time and space complexities of evaluating queries

Relational Algebra (Motivation)

- **SQL¹ and RA express the same queries:**
 - Each SQL query can be expressed by an RA expression
 - Each RA expression can be expressed by a SQL query.
- **Query optimization:** Rewrite rules can transform a RA expression into another equivalent RA expression that is more efficient to evaluate
- **Query evaluation:** Different algorithms and data structures can be associated with RA operations for efficient evaluations

¹without aggregation functions

Syntax of RA

- The relational algebra is a **typed** language of expressions
- Starting from relations and constants, RA expressions are inductively built using the operators

Operator	Algebraic notation
cartesian product	\times
selection	$\sigma(\cdot)$
projection	$\pi_{\dots}(\cdot)$
union	\cup
intersection	\cap
difference	$-$

- Each RA expression has a **schema** which is its type
- We will use $E(A_1, \dots, A_m)$ to denote a RA expression E with attribute schema (A_1, \dots, A_m)

RA basic expressions: relations

- Let $R(A_1, \dots, A_m)$ be a relation
- Then R is a RA expression with schema (A_1, \dots, A_m)
- Its value is the set of tuples from R , i.e., the relation instance associated with R
- R can be expressed by SQL query

```
SELECT   $r.A_1, \dots, r.A_m$ 
FROM     $R$   $r$ 
```

RA basic expressions: constants

- Let A be an attribute and let \mathbf{a} be a constant
- Then $(A:\mathbf{a})$ is a RA expression with schema (A)
- Its value is a unary relation containing the single tuple (\mathbf{a})
- $(A : \mathbf{a})$ can be expressed by SQL query

SELECT \mathbf{a} AS A

A

\mathbf{a}

RA operators: the cartesian product \times

- Let $R(A_1, \dots, A_m)$ and $S(B_1, \dots, B_n)$ be two relations with **non-overlapping** schemas
- Then $R \times S$ is a RA expression with schema $(A_1, \dots, A_m, B_1, \dots, B_n)$
- Its value is the set of all tuples that can be obtained by **pairing** *each* tuple r of R with *each* tuple s of S
- $R \times S$ can be expressed by the SQL query

```
SELECT   $r.A_1, \dots, r.A_m, s.B_1, \dots, s.B_n$   
FROM     $R r, S s$ 
```

- Observe that $|R \times S| = |R| |S|$

RA operators: the cartesian product \times (example)

R

A	B	C
a_1	b_1	c_1
a_2	b_2	c_2
a_3	b_3	c_3

S

D
1
2

$R \times S$

A	B	C	D
a_1	b_1	c_1	1
a_2	b_2	c_2	1
a_3	b_3	c_3	1
a_1	b_1	c_1	2
a_2	b_2	c_2	2
a_3	b_3	c_3	2

RA operators: the cartesian product \times (**building search space**)

- $R \times S$ provides in **single** relation all the information that is present in R and S
- $R \times S$ can be interpreted as the **search space**
- It provides a mechanism to associate each tuple r of R with each tuple s of S into a tuple (r, s) whose components can be compared
 - I.e., the components $r.A_1, \dots, r.A_m$ of r can be compared with the components $s.B_1, \dots, s.B_n$ of s

RA operators: selection $\sigma_{A \theta a}$ (Introduction)

- Let $R(A_1, \dots, A_m)$ be a relation
- Let a be a constant value
- Let θ be one of the comparison operators $=, \neq, <, \leq, >, \geq$
- Then $\sigma_{A_i \theta a}(R)$ is a RA expression with schema (A_1, \dots, A_m)
- Its value consists of the tuples r in R such that $r.A_i \theta a$ is true
- $\sigma_{A_i \theta a}(R)$ is expressed in SQL by

```
SELECT   $r.A_1, \dots, r.A_m$ 
FROM     $R$ 
WHERE    $r.A_i \theta a$ 
```

- Observe that $|\sigma_{A_i \theta a}(R)| \leq |R|$

RA operators: selection $\sigma_{A\theta a}$ (Example)

R

A	B	C
a_1	b_1	c_1
a_2	b_2	c_2
a_2	b_1	c_2
a_3	b_3	c_2

$\sigma_{B=b_1}(R)$

A	B	C
a_1	b_1	c_1
a_2	b_1	c_2

$\sigma_{A \neq a_3}(R)$

A	B	C
a_1	b_1	c_1
a_2	b_2	c_2
a_2	b_1	c_2

RA operators: selection $\sigma_{A \theta B}$ (Introduction)

- Let $R(A_1, \dots, A_m)$ be a relation
- Let θ be one of the comparison operators $=, \neq, <, \leq, >, \geq$
- Then $\sigma_{A_i \theta A_j}(R)$ is a RA expression with schema (A_1, \dots, A_m)
- Its value is the set of tuples r of R such that $r.A_i \theta r.A_j$ is true
- $\sigma_{A_i \theta A_j}(R)$ is expressed in SQL by

```
SELECT   $r.A_1, \dots, r.A_m$ 
FROM     $R$ 
WHERE    $r.A_i \theta r.A_j$ 
```

- Observe that $|\sigma_{A_i \theta A_j}(R)| \leq |R|$

RA operators: selection $\sigma_{A=B}$ (Example)

R

A	B	C
a_1	b_1	a_1
a_2	b_2	c_2
a_2	b_3	a_2

$\sigma_{A=C}(R)$

A	B	C
a_1	b_1	a_1
a_2	b_3	a_2

$\sigma_{A \neq C}(R)$

A	B	C
a_2	b_2	c_2

RA operators: selection σ

- A selection acts as a **filter** on R
- It provides a **horizontal slice** of R

RA operators: projection π (Introduction)

- Let $R(A_1, \dots, A_m)$ be a relation and let (B_1, \dots, B_k) be a non-empty list of k attributes of R
- Then $\pi_{B_1, \dots, B_k}(R)$ is a RA expression with schema (B_1, \dots, B_k)
- Its value is the relation $\{(r.B_1, \dots, r.B_k) \mid r \in R\}$
- In other words, each tuple r of R is "projected" on its (B_1, \dots, B_k) components.
- $\pi_{B_1, \dots, B_k}(R)$ is expressed in SQL by

```
SELECT DISTINCT   $r.B_1, \dots, r.B_k$   
FROM              $R$   $r$ 
```

- Observe that $|\pi_{B_1, \dots, B_k}(R)| \leq |R|$

RA operators: projection π (Example)

R

A	B	C
a_1	b_1	c_1
a_2	b_2	c_2
a_1	b_3	c_1
a_2	b_2	c_1

$\pi_{A,C}(R)$

A	C
a_1	c_1
a_2	c_2
a_2	c_1

$\pi_{B,A}(R)$

B	A
b_1	a_1
b_2	a_2
b_3	a_1

Notice that duplicates are eliminated and that projection permits attributes to be permuted

RA operators: projection π

- A projection provides a mechanism to get certain columns from R
- It provides a **vertical slice** of R

Example query in RA and SQL

- Consider the relations Student(sid,sname,age) and Enroll(sno,cno)⁴
- "Find the sid and age of each enrolled student whose name is Ann."
- This query can be expressed by the RA expression

$$\pi_{\text{sid, age}}(\sigma_{\text{sid=sno}}(\sigma_{\text{sname='Ann'}}(\text{Student} \times \text{Enroll})))$$

- In SQL

```
SELECT  DISTINCT s.sid, s.age
FROM    Student s, Enroll e
WHERE   s.sname = 'Ann' AND s.sid = e.sno
```

⁴Notice that we have used two different names sid and sno to refer uniquely to a student. We did this to ensure that the relations have no commonly named attributes.

Boolean conditions in selection operations

- It is possible to extend the selection operations by permitting boolean combinations of basic conditions " $A_j \theta a$ " and " $A_i \theta A_j$ " using the boolean connectors \wedge , \vee , and \neg .
- "Find the sid and age of each enrolled student whose name is Ann."
- This query can then be expressed using the following RA expression

$$\pi_{\text{sid, age}}(\sigma_{\text{name} = \text{'Ann'} \wedge \text{sid} = \text{sno}}(\text{Student} \times \text{Enroll}))$$

- Notice that the structure of this query is nearly the same as that of its SQL equivalent

Some comments about attribute names (renaming)

- Cartesian products require that the attributes names of the participating expressions do not overlap
- There are various ways to overcome this by renaming attributes or using some conventions
- For example, let $R(A, B)$ and $S(B, C)$ be two relations. We can not write $R \times S$ since attribute B occurs in the schemas of R and S . Nonetheless we will permit this and agree that the output schema of $R \times S$ is $(A, R.B, S.B, C)$ where we have used the relation names R and S to differentiate the B attribute in R from the B attribute in S

Some comments about attribute names

- We can also not write $R \times R$. To overcome this problem, we assume that for each relation name R , there is a series of relations R_1, R_2, \dots that are "copies" of R .
- Then instead of writing $R \times R$, we write $R_1 \times R_2$ and the output schema will be $(R_1.A, R_1.B, R_2.A, R_2.B)$
- In SQL, we can always use the **AS** clause to deal with attribute renaming. We will not give the details of attribute renaming in general but assume that it can always be done appropriately.

Some comments about attribute names (Example)

- "Find the sids of students who take at least two courses."

$$\pi_{E_1.sid}(\sigma_{E_1.sid = E_2.sid \wedge E_1.cno \neq E_2.cno}(Enroll_1 \times Enroll_2))$$

where E_1 and E_2 are abbreviations for $Enroll_1$ and $Enroll_2$, respectively

Boolean set operations (union, intersection, and set difference)

- Let E_1 and E_2 be two RA expressions with the **same** schema then

- $E_1 \cup E_2$,
- $E_1 \cap E_2$, and
- $E_1 - E_2$

are RA expressions with the same schema

- If Q_{E_1} and Q_{E_2} are the SQL queries corresponding to E_1 and E_2 , respectively, then the expressions are expressed in SQL by
 - $Q_{E_1} \text{ UNION } Q_{E_2}$,
 - $Q_{E_1} \text{ INTERSECT } Q_{E_2}$, and
 - $Q_{E_1} \text{ EXCEPT } Q_{E_2}$

Examples

- In this example, assume relation *Enroll*(*sid*, *cno*, *grade*)
- "Find the sid of each student who received an 'A' in some course."

$$\pi_{sid}(\sigma_{grade='A'}(Enroll))$$

- "Find the sid of each student who received an 'A' in some course and a 'B' in some course."

$$\pi_{sid}(\sigma_{grade='A'}(Enroll)) \cap \pi_{sid}(\sigma_{grade='B'}(Enroll))$$

- "Find the sid of each student who did not receive an 'A' in any course."

$$\pi_{sid}(Student) - \pi_{sid}(\sigma_{grade='A'}(Enroll))$$

- "Find the sids of students who are enrolled in at least one course"

$$\pi_{sid}(Enroll)$$

- "Find the sids of students who take at least two courses"

$$\pi_{E1.sid}(\sigma_{E1.sid = E2.sid \wedge E1.cno \neq E2.cno}(Enroll_1 \times Enroll_2))$$

- "Find the sids of students who take exactly one course"

$$\pi_{sid}(Enroll) - \pi_{E1.sid}(\sigma_{E1.sid = E2.sid \wedge E1.cno \neq E2.cno}(Enroll_1 \times Enroll_2))$$

Examples (**NOT ONLY** and **ONLY**)

- "Find the cnos of CS courses."

$$\pi_{cno}(\sigma_{dept='CS'}(Course))$$

Denote this expression by **CS**

- "Find sids of students who do **not only** take CS courses."

$$\pi_{sid}(Enroll - (\pi_{sid}(Student) \times CS))$$

- "Find sids of students who **only** take CS courses."

$$\pi_{sid}(Student) - \pi_{sid}(Enroll - (\pi_{sid}(Student) \times CS))$$

Examples (**NOT ALL** and **ALL**)

- "Find sids of students who do **not** take **all** CS courses"

$$\pi_{sid}((\pi_{sid}(Student) \times CS) - Enroll)$$

- "Find sids of students who take **all** CS courses"

$$\pi_{sid}(Student) - \pi_{sid}(\pi_{sid}(Student) \times CS) - Enroll)$$

Examples

Consider the relation $Person(pid, age)$

- "Find the pids of persons who are not the youngest."

$$\pi_{P_1.pid}(\sigma_{P_1.age > P_2.age}(Person_1 \times Person_2))$$

- "Find the pids of the youngest persons."

$$\pi_{pid}(Person) - \pi_{P_1.pid}(\sigma_{P_1.age > P_2.age}(Person_1 \times Person_2))$$