

CompSci 516 Database Systems

Lecture 4 Relational Algebra and Relational Calculus

Instructor: Sudeepa Roy

Duke CS, Spring 2022

CompSci 516: Database Systems

1

1

Announcements

- In-person classes starting Thursday
 - Also live streaming and recording

Duke CS, Spring 2022

CompSci 516: Database Systems

2

2

Today's topics

- Relational Algebra (RA) and Relational Calculus (RC)
- Reading material
 - [RG] Chapter 4 (RA, RC)
 - [GUW] Chapters 2.4, 5.1, 5.2

Acknowledgement:
The following slides have been created adapting the instructor material of the [RG] book provided by the authors Dr. Ramakrishnan and Dr. Gehrke.

Duke CS, Spring 2022

CompSci 516: Database Systems

3

3

Relational Query Languages

Duke CS, Spring 2022

CompSci 516: Database Systems

4

4

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 intended to be used for complex calculations
 - QLS support easy, efficient access to large data sets

Duke CS, Spring 2022

CompSci 516: Database Systems

5

5

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 (**Non-operational, declarative, or procedural**)
- **Note: Declarative (RC, SQL) vs. Operational (RA)**

Duke CS, Spring 2022

CompSci 516: Database Systems

6

6

Preliminaries (recap)

- 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**
 - query will run regardless of instance
 - 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

Duke CS, Spring 2022

CompSci 516: Database Systems

7

7

Example Schema and Instances

Sailors(sid, sname, rating, age)

Boats(bid, bname, color)

Reserves(sid, bid, day)

S1

sid	sname	rating	age
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

R1

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

8

Logic Notations

- \exists There exists
- \forall For all
- \wedge Logical AND
- \vee Logical OR
- \neg NOT
- \Rightarrow Implies

9

11

Relational Algebra (RA)

10

Relational Algebra

- Takes one or more relations as input, and produces a relation as output
 - operator
 - operand
 - semantic
 - so an algebra!
- Since each operation returns a relation, **operations can be composed**
 - Algebra is “closed”

Duke CS, Spring 2022

CompSci 516: Database Systems

11

Relational Algebra

- **Basic operations:**
 - **Selection (σ)** Selects a subset of rows from relation
 - **Projection (π)** Deletes unwanted columns from relation.
 - **Cross-product (\times)** Allows us to combine two relations.
 - **Set-difference ($-$)** Tuples in reln. 1, but not in reln. 2.
 - **Union (\cup)** Tuples in reln. 1 or in reln. 2.
- **Additional operations:**
 - **Intersection (\cap)**
 - **join (\bowtie)**
 - **division ($/$)**
 - **renaming (ρ)**
 - Not essential, but (very) useful.

Duke CS, Spring 2022

CompSci 516: Database Systems

12

12

Projection

S_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

- Deletes attributes that are not in projection list.
- 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 (performance)

sname	rating
yuppy	9
lubber	8
guppy	5
rusty	10

$\pi_{sname, rating}(S_2)$

age
35.0
55.5

$\pi_{age}(S_2)$

Duke CS, Fall 2016 CompSci 516: Data Intensive Computing Systems 10

13

Selection

S_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

- Selects rows that satisfy **selection condition**
- No duplicates in result. Why?
- Schema of result identical to schema of (only) input relation

sid	sname	rating	age
28	yuppy	9	35.0
58	rusty	10	35.0

$\sigma_{rating > 8}(S_2)$

Duke CS, Fall 2016 CompSci 516: Data Intensive Computing Systems 11

14

Composition of Operators

- Result relation can be the input for another relational algebra operation
 - Operator composition

sid	sname	rating	age
28	yuppy	9	35.0
58	rusty	10	35.0

$\sigma_{rating > 8}(S_2)$

sname	rating
yuppy	9
rusty	10

$\pi_{sname, rating}(\sigma_{rating > 8}(S_2))$

Duke CS, Spring 2016 CompSci 516: Data Intensive Computing Systems 12

15

Union, Intersection, Set-Difference

S_1

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

S_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

- All of these operations take two input relations, which must be **union-compatible**:
 - Same number of fields.
 - 'Corresponding' fields have the same type
 - same schema as the inputs

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

$S_1 \cup S_2$

Duke CS, Spring 2016 CompSci 516: Data Intensive Computing Systems 12

16

Union, Intersection, Set-Difference

S_1

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

S_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

- Note: no duplicate
 - "Set semantic"
 - SQL: **UNION**
 - SQL allows "bag semantic" as well: **UNION ALL**

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

$S_1 \cup S_2$

Duke CS, Spring 2016 CompSci 516: Data Intensive Computing Systems 12

17

Union, Intersection, Set-Difference

S_1

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

S_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

$S_1 - S_2$

$S_1 \cap S_2$

Duke CS, Spring 2016 CompSci 516: Data Intensive Computing Systems 13

18

Cross-Product

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

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

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

(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

Duke CS, Spring 2022 CompSci 516: Database Systems 19

19

Renaming Operator ρ

$(\rho_{sid \rightarrow sid1} S1) \times (\rho_{sid \rightarrow sid1} R1)$

or

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

C is the new relation name

(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

■ In general, can use $\rho(\langle Temp \rangle, \langle RA-expression \rangle)$

Duke CS, Spring 2022 CompSci 516: Database Systems 20

20

Joins

$R \bowtie_c S = \sigma_c(R \times S)$

(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

Duke CS, Spring 2022 CompSci 516: Database Systems 21

21

Find names of sailors who've reserved boat #103

Sailors(sid, sname, rating, age)
 Boats(bid, bname, color)
 Reserves(sid, bid, day)

Duke CS, Spring 2022 CompSci 516: Database Systems 22

22

Find names of sailors who've reserved boat #103

Sailors(sid, sname, rating, age)
 Boats(bid, bname, color)
 Reserves(sid, bid, day)

- Solution 1: $\pi_{sname}((\sigma_{bid=103} Reserves) \bowtie Sailors)$
- Solution 2: $\pi_{sname}(\sigma_{bid=103}(Reserves \bowtie Sailors))$

Duke CS, Spring 2022 CompSci 516: Database Systems 23

23

Expressing an RA expression as a Tree

Sailors(sid, sname, rating, age)
 Boats(bid, bname, color)
 Reserves(sid, bid, day)

Also called a **logical query plan**

$\pi_{sname}((\sigma_{bid=103} Reserves) \bowtie Sailors)$

Duke CS, Spring 2022 CompSci 516: Database Systems 24

24

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

Sailors(sid, sname, rating, age)
 Boats(bid, bname, color)
 Reserves(sid, bid, day)

Use of rename operation

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

Can also define Tempboats using union
 Try the "AND" version yourself

Duke CS, Spring 2022 CompSci 516: Database Systems 25

25

What about aggregates?

Sailors(sid, sname, rating, age)
 Boats(bid, bname, color)
 Reserves(sid, bid, day)

- Extended relational algebra
- $\gamma_{age, avg(rating)} \rightarrow avgr$ Sailors
- Also extended to "bag semantic": allow duplicates
 - Take into account cardinality
 - R and S have tuple t resp. m and n times
 - $R \cup S$ has t m+n times
 - $R \cap S$ has t min(m, n) times
 - $R - S$ has t max(0, m-n) times
 - sorting(τ), duplicate removal (δ) operators

Duke CS, Spring 2022 CompSci 516: Database Systems 26

26

Relational Calculus (RC)

Duke CS, Spring 2022 CompSci 516: Database Systems 27

27

Relational Calculus

- RA is procedural
 - $\pi_A(\sigma_{A=a} R)$ and $\sigma_{A=a}(\pi_A R)$ are equivalent but different expressions
- RC
 - non-procedural and declarative
 - describes a set of answers without being explicit about how they should be computed
- TRC (tuple relational calculus)
 - variables take tuples as values
 - we will primarily do TRC
- DRC (domain relational calculus)
 - variables range over field values

Duke CS, Spring 2022 CompSci 516: Database Systems 28

28

TRC: example

Sailors(sid, sname, rating, age)
 Boats(bid, bname, color)
 Reserves(sid, bid, day)

- Find the name and age of all sailors with a rating above 7

\exists There exists

$\{P \mid \exists S \in \text{Sailors} (S.\text{rating} > 7 \wedge P.\text{sname} = S.\text{sname} \wedge P.\text{age} = S.\text{age})\}$

- P is a tuple variable
 - with exactly two fields sname and age (schema of the output relation)
 - $P.\text{sname} = S.\text{sname} \wedge P.\text{age} = S.\text{age}$ gives values to the fields of an answer tuple
- Use parentheses, $\forall \exists \vee \wedge > < = \neq \sim$ etc as necessary
- $A \Rightarrow B$ is very useful too

next slide

Duke CS, Spring 2022 CompSci 516: Database Systems 29

29

$A \Rightarrow B$

- A "implies" B
- Equivalently, if A is true, B must be true
- Equivalently, $\neg A \vee B$, i.e.
 - either A is false (then B can be anything)
 - otherwise (i.e. A is true) B must be true

Duke CS, Spring 2022 CompSci 516: Database Systems 30

30

Useful Logical Equivalences

- $\forall x P(x) = \neg \exists x [\neg P(x)]$

\exists	There exists
\forall	For all
\wedge	Logical AND
\vee	Logical OR
\neg	NOT

- $\neg(P \vee Q) = \neg P \wedge \neg Q$
- $\neg(P \wedge Q) = \neg P \vee \neg Q$

} de Morgan's laws

– Similarly, $\neg(\neg P \vee Q) = P \wedge \neg Q$ etc.

- $A \Rightarrow B = \neg A \vee B$

Duke CS, Spring 2022 CompSci 516: Database Systems 31

31

TRC: example

Sailors(sid, sname, rating, age)
Boats(bid, bname, color)
Reserves(sid, bid, day)

- Find the names of sailors who have reserved at least two boats

Duke CS, Spring 2022 CompSci 516: Database Systems 32

32

TRC: example

Sailors(sid, sname, rating, age)
Boats(bid, bname, color)
Reserves(sid, bid, day)

- Find the names of sailors who have reserved at least two boats

{P | $\exists S \in \text{Sailors} (\exists R1 \in \text{Reserves} \exists R2 \in \text{Reserves} (S.\text{sid} = R1.\text{sid} \wedge S.\text{sid} = R2.\text{sid} \wedge R1.\text{bid} \neq R2.\text{bid}) \wedge P.\text{sname} = S.\text{sname})$ }

Duke CS, Spring 2022 CompSci 516: Database Systems 33

33

TRC: example

Sailors(sid, sname, rating, age)
Boats(bid, bname, color)
Reserves(sid, bid, day)

- Find the names of sailors who have reserved all boats
- Called the "Division" operation in RA

Duke CS, Spring 2022 CompSci 516: Database Systems 34

34

TRC: example

Sailors(sid, sname, rating, age)
Boats(bid, bname, color)
Reserves(sid, bid, day)

- Find the names of sailors who have reserved all boats

{P | $\exists S \in \text{Sailors} [\forall B \in \text{Boats} (\exists R \in \text{Reserves} (S.\text{sid} = R.\text{sid} \wedge R.\text{bid} = B.\text{bid}))] \wedge (P.\text{sname} = S.\text{sname})$ }

Duke CS, Spring 2022 CompSci 516: Database Systems 35

35

TRC: example

Sailors(sid, sname, rating, age)
Boats(bid, bname, color)
Reserves(sid, bid, day)

- Find the names of sailors who have reserved all red boats

How will you change the previous TRC expression?

Duke CS, Spring 2022 CompSci 516: Database Systems 36

36

TRC: example

Sailors(sid, sname, rating, age)
 Boats(bid, bname, color)
 Reserves(sid, bid, day)

- Find the names of sailors who have reserved all red boats

Recall that $A \Rightarrow B$ is logically equivalent to $\neg A \vee B$
 so \Rightarrow can be avoided, but it is cleaner and more intuitive

Duke CS, Spring 2022 CompSci 516: Database Systems 37

37

DRC: example

Sailors(sid, sname, rating, age)
 Boats(bid, bname, color)
 Reserves(sid, bid, day)

- Find the name and age of all sailors with a rating above 7

TRC:
 $\{P \mid \exists S \in \text{Sailors} (S.\text{rating} > 7 \wedge P.\text{name} = S.\text{name} \wedge P.\text{age} = S.\text{age})\}$

DRC:
 $\{<N, A> \mid \exists <I, N, T, A> \in \text{Sailors} \wedge T > 7\}$

- Variables are now domain variables
- We will use TRC
 - both are equivalent
- Another option to write coming soon!

Duke CS, Spring 2022 CompSci 516: Database Systems 38

38

The famous "Beers" database

Drinkers **Frequent** Bars
 "X" times a week

Bars
 Each has an address

Bars **Serve** Beers
 At price "Y"

Drinkers **Likes** Beers

Beers
 Each has a brewer

39

39

See online database for more tuples

"Beers" as a Relational Database

name	address
The Edge	108 Morris Street
Satisfaction	905 W. Main Street

bar	beer	price
The Edge	Budweiser	2.50
The Edge	Corona	3.00
Satisfaction	Budweiser	2.25

Name	brewer
Budweiser	Anheuser-Busch Inc.
Corona	Grupo Modelo
Dixie	Dixie Brewing

drinker	bar	times_a_week
Ben	Satisfaction	2
Dan	The Edge	1
Dan	Satisfaction	2

name	address
Amy	100 W. Main Street
Ben	101 W. Main Street
Dan	300 N. Duke Street

drinker	beer
Amy	Corona
Dan	Budweiser
Dan	Corona
Ben	Budweiser

40

40

More Examples: RC

UNDERSTAND THE DIFFERENCE IN ANSWERS
 FOR ALL FOUR DRINKERS

Acknowledgement: examples and slides by Profs. Balazinska and Suciu, and the [GUW] book

Duke CS, Spring 2022 CompSci 516: Database Systems 41

41

Likes(drinker, beer)
 Frequent(drinker, bar)
 Serves(bar, beer)

Drinker Category 1

Find drinkers that frequent some bar that serves some beer they like.

...

42

42

Likes(drinker, beer)
 Frequent(drinker, bar)
 Serves(bar, beer)

Drinker Category 1

Find drinkers that frequent some bar that serves some beer they like.

43

43

Likes(drinker, beer)
 Frequent(drinker, bar)
 Serves(bar, beer)

Drinker Category 2

Find drinkers that frequent some bar that serves some beer they like.

Find drinkers that frequent only bars that serves some beer they like.

...

Free HW question hint!

44

44

Likes(drinker, beer)
 Frequent(drinker, bar)
 Serves(bar, beer)

Drinker Category 2

Find drinkers that frequent some bar that serves some beer they like.

Find drinkers that frequent only bars that serve some beer they like.

45

45

Likes(drinker, beer)
 Frequent(drinker, bar)
 Serves(bar, beer)

Drinker Category 3

Find drinkers that frequent some bar that serves some beer they like.

Find drinkers that frequent only bars that serve some beer they like.

Find drinkers that frequent some bar that serves only beers they like.

...

46

46

Likes(drinker, beer)
 Frequent(drinker, bar)
 Serves(bar, beer)

Drinker Category 3

Find drinkers that frequent some bar that serves some beer they like.

Find drinkers that frequent only bars that serve some beer they like.

Find drinkers that frequent some bar that serves only beers they like.

47

47

Likes(drinker, beer)
 Frequent(drinker, bar)
 Serves(bar, beer)

Drinker Category 4

Find drinkers that frequent some bar that serves some beer they like.

Find drinkers that frequent only bars that serve some beer they like.

Find drinkers that frequent some bar that serves only beers they like.

Find drinkers that frequent only bars that serve only beer they like.

...

48

48

Likes(drinker, beer)
Frequents(drinker, bar)
Serves(bar, beer)

Drinker Category 4

Find drinkers that frequent some bar that serves some beer they like.

Find drinkers that frequent only bars that serve some beer they like.

Find drinkers that frequent some bar that serves only beers they like.

Find drinkers that frequent only bars that serve only beer they like.

49

49

Why should we care about RC

- RC is declarative, like SQL, and unlike RA (which is operational)
- Gives foundation of database queries in first-order logic
 - you cannot express all aggregates in RC, e.g. cardinality of a relation or sum (possible in extended RA and SQL)
 - still can express conditions like “at least two tuples” (or any constant)
- RC expression may be much simpler than SQL queries
 - and easier to check for correctness than SQL
 - power to use \forall and \Rightarrow
 - then you can systematically go to a “correct” SQL or RA query

50

50

Likes(drinker, beer)
Frequents(drinker, bar)
Serves(bar, beer)

Drinker category 5!

From RC to SQL

Query: Find drinkers that like some beer (so much) that they frequent all bars that serve it

$$\{x \mid \exists L \in \text{Likes} (L.\text{drinker} = x.\text{drinker}) \wedge [\forall S \in \text{Serves} (L.\text{beer} = S.\text{beer}) \Rightarrow \exists F \in \text{Frequents} [(F.\text{drinker} = L.\text{drinker}) \wedge (F.\text{bar} = S.\text{bar})]]\}$$

51

51

Likes(drinker, beer)
Frequents(drinker, bar)
Serves(bar, beer)

From RC to SQL (or RA)

Query: Find drinkers that like some beer so much that they frequent all bars that serve it

$$\{x \mid \exists L \in \text{Likes} (L.\text{drinker} = x.\text{drinker}) \wedge [\forall S \in \text{Serves} [(L.\text{beer} = S.\text{beer}) \Rightarrow \exists F \in \text{Frequents} [(F.\text{drinker} = L.\text{drinker}) \wedge (F.\text{bar} = S.\text{bar})]]]\}$$

$$\equiv \{x \mid \exists L \in \text{Likes} (L.\text{drinker} = x.\text{drinker}) \wedge [\forall S \in \text{Serves} [\neg (L.\text{beer} = S.\text{beer}) \vee \exists F \in \text{Frequents} [(F.\text{drinker} = L.\text{drinker}) \wedge (F.\text{bar} = S.\text{bar})]]]\}$$

Step 1: Replace \forall with \exists using de Morgan's Laws

$$Q(x) = \exists y. \text{Likes}(x, y) \wedge [\neg \exists S \in \text{Serves} [(L.\text{beer} = S.\text{beer}) \wedge \neg \exists F \in \text{Frequents} [(F.\text{drinker} = L.\text{drinker}) \wedge (F.\text{bar} = S.\text{bar})]]]$$

SQL or RA does not have \forall !
Now you got all \exists and \neg expressible in RA/SQL

52

52

Likes(drinker, beer)
Frequents(drinker, bar)
Serves(bar, beer)

From RC to SQL

Query: Find drinkers that like some beer so much that they frequent all bars that serve it

$$\exists L \in \text{Likes} \wedge \neg \exists S \in \text{Serves} [(L.\text{beer} = S.\text{beer}) \wedge \neg \exists F \in \text{Frequents} [(F.\text{drinker} = L.\text{drinker}) \wedge (F.\text{bar} = S.\text{bar})]]$$

Step 2: Translate into SQL

```
SELECT DISTINCT L.drinker
FROM Likes L
WHERE not exists
(SELECT S.bar
FROM Serves S
WHERE L.beer=S.beer
AND not exists (SELECT *
FROM Frequents F
WHERE F.drinker=L.drinker
AND F.bar=S.bar))
```

We will see a “methodical and correct” translation through “safe queries” in Datalog

53

53

Summary

- You learnt three query languages for the Relational DB model
 - SQL
 - RA
 - RC
- All have their own purposes
- You should be able to write a query in all three languages and convert from one to another
 - However, you have to be careful, not all “valid” expressions in one may be expressed in another
 - $\{S \mid \neg (S \in \text{Sailors})\}$ – infinitely many tuples – an “unsafe” query
 - More when we do “Datalog”, also see Ch. 4.4 in [RG]

54

54