

Database Systems CSE 414

Lectures 8: Relational Algebra (Ch. 2.4, & 5.1)

CSE 414 - Fall 2017

1

Announcements

- WQ4 is posted and due on Nov. 3, 11pm
- HW2 will be due next Monday 11pm
- **Log into Azure web site using your outlook.com email address today if you have not done so already**
 - Otherwise the TA cannot give you the Azure access code

CSE 414 - Fall 2017

2

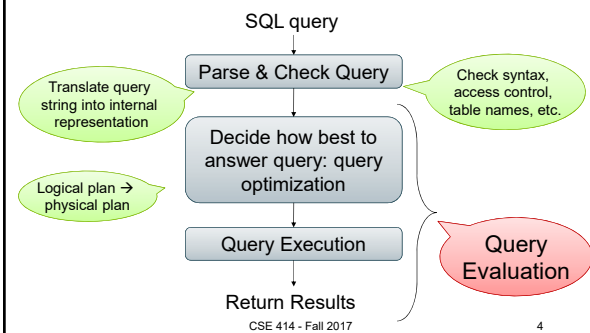
Where We Are

- Motivation for using a DBMS for managing data
- SQL:
 - Declaring the schema for our data (CREATE TABLE)
 - Inserting data one row at a time or in bulk (INSERT/import)
 - Modifying the schema and updating the data (ALTER/UPDATE)
 - Querying the data (SELECT)
- **Next step: More knowledge of how DBMSs work**
 - Client-server architecture
 - Relational algebra and query execution

CSE 414 - Fall 2017

3

Query Evaluation Steps



CSE 414 - Fall 2017

4

The WHAT and the HOW

- SQL = **WHAT** we want to get from the data
- Relational Algebra = **HOW** to get the data we want
- Move from **WHAT** to **HOW** is **query optimization**
 - SQL ~> Relational Algebra ~> Physical Plan
 - Relational Algebra = Logical Plan

CSE 414 - Fall 2017

5

Relational Algebra

CSE 414 - Fall 2017

6

Sets vs. Bags

- Sets: {a,b,c}, {a,d,e,f}, {}, . . .
- Bags: {a, a, b, c}, {b, b, b, b}, . . .

Relational Algebra has two semantics:

- Set semantics = standard Relational Algebra
- Bag semantics = extended Relational Algebra

DB systems implement bag semantics (Why?)

CSE 414 - Fall 2017

7

Relational Algebra Operators

- Union \cup , intersection \cap , difference -
- Selection σ
- Projection π (Π)
- Cartesian product \times , join \bowtie
- Rename ρ
- Duplicate elimination δ
- Grouping and aggregation γ
- Sorting τ

RA

Extended RA

CSE 414 - Fall 2017

8

Union and Difference

$$\begin{array}{l} R1 \cup R2 \\ R1 - R2 \end{array}$$

What do they mean over bags ?

CSE 414 - Fall 2017

9

What about Intersection ?

- Derived operator using minus
- Derived using join (will explain later)

$$R1 \cap R2 = R1 - (R1 - R2)$$

$$R1 \cap R2 = R1 \bowtie R2$$

CSE 414 - Fall 2017

10

Selection

- Returns all tuples that satisfy a condition

$$\sigma_c(R)$$

- Examples
 - $\sigma_{\text{Salary} > 40000}$ (Employee)
 - $\sigma_{\text{name} = \text{"Smith"}}$ (Employee)
- The condition c can be $=, <, \leq, >, \geq, \neq$ combined with AND, OR, NOT

CSE 414 - Fall 2017

11

Employee

SSN	Name	Salary
1234545	John	20000
5423341	Smith	60000
4352342	Fred	50000

$\sigma_{\text{Salary} > 40000}$ (Employee)

SSN	Name	Salary
5423341	Smith	60000
4352342	Fred	50000

CSE 414 - Fall 2017

12

Projection

- Eliminates columns

$$\pi_{A_1, \dots, A_n}(R)$$

- Example: project social-security number and names:

- $\Pi_{SSN, Name}(Employee)$
- Answer(SSN, Name)

Different semantics over sets or bags! Why?

CSE 414 - Fall 2017

13

Employee

SSN	Name	Salary
1234545	John	20000
5423341	John	60000
4352342	John	20000

$\pi_{Name, Salary}(Employee)$

Name	Salary
John	20000
John	60000
John	20000

Name	Salary
John	20000
John	60000

Bag semantics

Set semantics

Which is more efficient?

14

Composing RA Operators

Patient

no	name	zip	disease
1	p1	98125	flu
2	p2	98125	heart
3	p3	98120	lung
4	p4	98120	heart

$\pi_{zip, disease}(Patient)$

zip	disease
98125	flu
98125	heart
98120	lung
98120	heart

$\sigma_{disease='heart'}(Patient)$

no	name	zip	disease
2	p2	98125	heart
4	p4	98120	heart

$\pi_{zip, disease}(\sigma_{disease='heart'}(Patient))$

zip	disease
98125	heart
98120	heart

CSE 414 - Fall 2017

15

Cartesian Product

- Each tuple in R1 with each tuple in R2

$$R1 \times R2$$

- Rare in practice; mainly used to express joins

CSE 414 - Fall 2017

16

Cross-Product Example

Employee

Name	SSN
John	999999999
Tony	777777777

Dependent

EmpSSN	DepName
999999999	Emily
777777777	Joe

Employee \times Dependent

Name	SSN	EmpSSN	DepName
John	999999999	999999999	Emily
John	999999999	777777777	Joe
Tony	777777777	999999999	Emily
Tony	777777777	777777777	Joe

CSE 414 - Fall 2017

17

Renaming

- Changes the schema, not the instance

$$\rho_{B_1, \dots, B_n}(R)$$

- Example:

- $\rho_{N, S}(Employee) \rightarrow Answer(N, S)$

Not really used by systems, but needed on paper

CSE 414 - Fall 2017

18

Natural Join

$$R1 \bowtie R2$$

- Meaning: $R1 \bowtie R2 = \pi_A(\sigma_\theta(R1 \times R2))$
- Where:
 - Selection σ checks equality of **all common attributes** (attributes with same names)
 - Projection π eliminates duplicate **common attributes**

CSE 414 - Fall 2017

19

Natural Join Example

R	A	B	S	B	C
	X	Y		Z	U
	X	Z		V	W
	Y	Z		Z	V
	Z	V			

$R \bowtie S = \pi_{A,B,C}(\sigma_{R.B=S.B}(R \times S))$

A	B	C
X	Z	U
X	Z	V
Y	Z	U
Y	Z	V
Z	V	W

CSE 414 - Fall 2017

20

Natural Join Example 2

AnonPatient P

age	zip	disease
54	98125	heart
20	98120	flu

Voters V

name	age	zip
p1	54	98125
p2	20	98120

$P \bowtie V$

age	zip	disease	name
54	98125	heart	p1
20	98120	flu	p2

CSE 414 - Fall 2017

21

Natural Join

- Given schemas $R(A, B, C, D)$, $S(A, C, E)$, what is the schema of $R \bowtie S$?
 - (A, B, C, D, E) through join on (A, C)
- Given $R(A, B, C)$, $S(D, E)$, what is $R \bowtie S$?
 - (A, B, C, D, E) through cross product
- Given $R(A, B)$, $S(A, B)$, what is $R \bowtie S$?
 - (A, B) through intersection

CSE 414 - Fall 2017

22

AnonPatient (age, zip, disease)
Voters (name, age, zip)

Theta Join

- A join that involves a predicate

$$R1 \bowtie_{\theta} R2 = \sigma_{\theta}(R1 \times R2)$$

- Here θ can be any condition
- For our voters/patients example:

$$P \bowtie_{P.zip = V.zip \text{ and } P.age \geq V.age - 1 \text{ and } P.age \leq V.age + 1} V$$

CSE 414 - Fall 2017

23

Equijoin

- A theta join where θ is an equality predicate
- By far the most used variant of join in practice

CSE 414 - Fall 2017

24

Equijoin Example

AnonPatient P

age	zip	disease
54	98125	heart
20	98120	flu

Voters V

name	age	zip
p1	54	98125
p2	20	98120

$$P \bowtie_{P.age=V.age} V$$

P.age	P.zip	P.disease	P.name	V.zip	V.age
54	98125	heart	p1	98125	54
20	98120	flu	p2	98120	20

Join Summary

- **Theta-join:** $R \bowtie_{\theta} S = \sigma_{\theta}(R \times S)$
 - Join of R and S with a join condition θ
 - Cross-product followed by selection θ
- **Equijoin:** $R \bowtie_{\theta} S = \sigma_{\theta}(R \times S)$
 - Join condition θ consists only of equalities
- **Natural join:** $R \bowtie S = \pi_A(\sigma_{\theta}(R \times S))$
 - Equijoin
 - Equality on **all** fields with same name in R and in S
 - Projection π_A drops all redundant attributes

So Which Join Is It ?

When we write $R \bowtie S$, we usually mean an equijoin, but we often omit the equality predicate when it is clear from the context

More Joins

- **Outer join**
 - Include tuples with no matches in the output
 - Use NULL values for missing attributes
 - Does not eliminate duplicate columns
- **Variants**
 - Left outer join
 - Right outer join
 - Full outer join

Outer Join Example

AnonPatient P

age	zip	disease
54	98125	heart
20	98120	flu
33	98120	lung

AnnonJob J

job	age	zip
lawyer	54	98125
cashier	20	98120

$$P \bowtie J$$

P.age	P.zip	disease	job	J.age	J.zip
54	98125	heart	lawyer	54	98125
20	98120	flu	cashier	20	98120
33	98120	lung	null	null	null

More Examples

```
Supplier(sno, sname, scity, sstate)
Part(pno, pname, psize, pcolor)
Supply(sno, pno, qty, price)
```

Name of supplier of parts with size greater than 10

$$\pi_{sname}(\text{Supplier} \bowtie \text{Supply} \bowtie (\sigma_{psize > 10}(\text{Part})))$$

Name of supplier of red parts or parts with size greater than 10

$$\pi_{sname}(\text{Supplier} \bowtie \text{Supply} \bowtie (\sigma_{psize > 10}(\text{Part}) \cup \sigma_{pcolor='red'}(\text{Part})))$$