

# Introduction to Data Management

## CSE 344

### Lectures 8: Relational Algebra

# Announcements

- Homework 3 is posted
  - Microsoft Azure Cloud services!
  - Use the promotion code you received
  - Due on 2/1
- Make sure you read the textbook!
  - Very good coverage of RA

# Where We Are

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

# Query Evaluation Steps

SQL query

Parse & Check Query

Translate query string into internal representation

Check syntax, access control, table names, etc.

Decide how best to answer query: query optimization

Logical plan → physical plan

Relational Algebra

Query Execution

Query Evaluation

Return Results

# The WHAT and the HOW

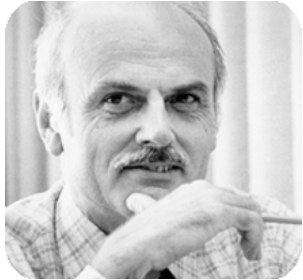
- SQL = **WHAT** we want to get from the data
- Relational Algebra = **HOW** to get the data we want
- The passage from **WHAT** to **HOW** is called **query optimization**
  - SQL → Logical Plan → Physical Plan
  - Logical plan expressed using relational algebra

# Relational Algebra

# Turing Awards in Data Management



Charles Bachman, 1973  
*IDS and CODASYL*



Ted Codd, 1981  
*Relational model*



Jim Gray, 1998  
*Transaction processing*



Michael Stonebraker, 2014  
*INGRES and Postgres*

# Sets v.s. Bags

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

Relational Algebra has two semantics:

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

DB systems implement bag semantics (Why?)



# Relational Algebra Operators

- Union  $\cup$ , intersection  $\cap$ , difference  $-$
- Selection  $\sigma$
- Projection  $\pi$
- Cartesian product  $\times$ , join  $\bowtie$
- Rename  $\rho$
- Duplicate elimination  $\delta$
- Grouping and aggregation  $\gamma$
- Sorting  $\tau$

RA

Extended RA

All operators take in 1 or more relations as inputs  
and return another relation

# Union and Difference

$$R1 \cup R2$$
$$R1 - R2$$

What do they mean over bags ?

# What about Intersection ?

- Derived operator using minus

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

– Only makes sense if result is  $\geq 0$

- Derived using join

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

– Only makes sense if R1 and R2 have the same schema

# Selection

- Returns all tuples which satisfy a condition

$$\sigma_c(R)$$

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

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

# Projection

- Eliminates columns

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

- Example: project social-security number and names:
  - $\pi_{\text{SSN}, \text{Name}}(\text{Employee}) \rightarrow \text{Answer}(\text{SSN}, \text{Name})$

Different semantics over sets or bags! Why?

Employee

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

$\Pi_{\text{Name,Salary}}$  (Employee)

Name	Salary
John	20000
John	60000
John	20000

Bag semantics

Name	Salary
John	20000
John	60000

Set semantics

Which is more efficient?

# Composing RA Operators

Patient

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

$\Pi_{\text{zip,disease}}(\text{Patient})$

zip	disease
98125	flu
98125	heart
98120	lung
98120	heart

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

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

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

zip	disease
98125	heart
98120	heart



# Cartesian Product

- Each tuple in R1 with each tuple in R2

$$R1 \times R2$$

- Rare in practice; mainly used to express joins

# Cross-Product Example

## Employee

Name	SSN
John	999999999
Tony	777777777

## Dependent

EmpSSN	DepName
999999999	Emily
777777777	Joe

## Employee X Dependent

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

# Renaming

- Changes the schema, not the instance

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

- Example:
  - Given Employee(Name, SSN)
  - $\rho_{N, S}(\text{Employee}) \rightarrow \text{Answer}(N, S)$

Not really used by systems, but needed on paper

# Natural Join

$$R1 \bowtie R2$$

- Meaning:  $R1 \bowtie R2 = \Pi_A(\sigma_\theta(R1 \times R2))$
- Where:
  - Selection  $\sigma_\theta$  checks equality of **all common attributes** (i.e., attributes with same names)
  - Projection  $\Pi_A$  eliminates duplicate **common attributes**

# Natural Join Example

**R**

A	B
X	Y
X	Z
Y	Z
Z	V

**S**

B	C
Z	U
V	W
Z	V

**R** ⋈ **S** =

$\Pi_{ABC}(\sigma_{R.B=S.B}(R \times S))$

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

# 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

join predicate:

$P.age = V.age$

AND

$P.zip = V.zip$

# Natural Join

- Given schemas  $R(A, B, C, D)$ ,  $S(A, C, E)$ , what is the schema of  $R \bowtie S$  ?
- Given  $R(A, B, C)$ ,  $S(D, E)$ , what is  $R \bowtie S$ ?
- Given  $R(A, B)$ ,  $S(A, B)$ , what is  $R \bowtie S$ ?

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  $\theta$  can be any condition
- No projection in this case!
- For our voters/patients example:

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



# Equijoin

- A theta join where  $\theta$  is an equality predicate
- Projection drops all redundant attributes

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

- By far the most used variant of join in practice
- What is the relationship with natural join?

# Equijoin Example

AnonPatient P

age	zip	disease
54	98125	heart
20	98120	flu

Voters V

name	age	zip
p1	54	98125
p2	20	98120

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

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

# Join Summary

- **Theta-join:**  $R \bowtie_{\theta} S = \sigma_{\theta} (R \times S)$ 
  - Join of R and S with a join condition  $\theta$
  - Cross-product followed by selection  $\theta$
- **Equijoin:**  $R \bowtie_{\theta} S = \pi_A (\sigma_{\theta} (R \times S))$ 
  - Join condition  $\theta$  consists only of equalities
  - Projection  $\pi_A$  drops all redundant attributes
- **Natural join:**  $R \bowtie S = \pi_A (\sigma_{\theta} (R \times S))$ 
  - Equality on **all** fields with same name in R and in S
  - Projection  $\pi_A$  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	<b>null</b>	33	98120

$\bowtie$  LOJ

$\ltimes$  ROJ

$\ltimes$  FOJ

# Some 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

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

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

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

Can be represented as trees as well (as seen from last class)