

# Introduction to Database Systems CSE 414

## Lecture 10: Relational Algebra

CSE 414 - Spring 2018

1

## Recap: Datalog

- Facts and Rules
- Selection, projection, join
- Recursive rules
- Grouping, aggregates
- Negation
- Safe vs unsafe rules
- Stratification

CSE 414 - Spring 2018

2

## Class Overview

- Unit 1: Intro
- Unit 2: Relational Data Models and Query Languages
  - Data models, SQL, Datalog, **Relational Algebra**
- Unit 3: Non-relational data
- Unit 4: RDBMS internals and query optimization
- Unit 5: Parallel query processing
- Unit 6: DBMS usability, conceptual design
- Unit 7: Transactions

CSE 414 - Spring 2018

3

## Relational Algebra

- Set-based algebra that manipulates relations
  - We will extend it to multisets / bags
- In SQL & Datalog we say *what* we want
- In RA we can express *how* to get it
- Every DBMS implementations converts a SQL query to RA in order to execute it
- An RA expression is called a *query plan*

CSE 414 - Spring 2018

4

## Why study yet another relational query language?

- RA is how SQL is implemented in DBMS
  - We will see more of this in a few weeks
- RA opens up opportunities for *query optimization*

CSE 414 - Spring 2018

5

## Basics

- Relations and attributes
- Functions that are applied to relations
  - Return relations
  - Can be composed together
  - Often displayed using a tree rather than linearly
  - Use Greek symbols:  $\sigma$ ,  $\pi$ ,  $\delta$ , etc

Relational algebra FTW!



CSE 414 - Spring 2018

6

## 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  $\tau$

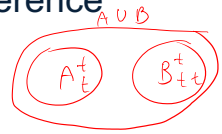
RA

Extended RA

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

## Union and Difference

$R1 \cup R2$   
 $R1 - R2$



Only make sense if  $R1, R2$  have the same schema

What do they mean over bags ?

CSE 414 - Spring 2018

8

## What about Intersection ?

- Derived operator using minus

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

- Derived using join (as we will see later)

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

CSE 414 - Spring 2018

9

## Selection

- Returns all tuples which satisfy a condition

$\sigma_c(R)$

- Examples

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

–  $\sigma_{\text{name} = \text{"Smith"}}(\text{Employee})$

- The condition  $c$  can be  $=, <, <=, >, >=, <>$  combined with AND, OR, NOT

CSE 414 - Spring 2018

10

Employee

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

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

SSN	Name	Salary
5423341	Smith	60000
4352342	Fred	50000

CSE 414 - Spring 2018

11

## Projection

- Eliminates columns

$\pi_{A1, \dots, An}(R)$

- 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

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

Name	Salary
John	20000
John	60000
John	20000

Name	Salary
John	20000
John	60000

Bag semantics                      Set semantics

Which is more efficient?

13

### 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 - Spring 2018                      14

### Cartesian Product

- Each tuple in R1 with each tuple in R2

$R1 \times R2$

- Rare in practice; mainly used to express joins

CSE 414 - Spring 2018                      15

### Cross-Product Example

Employee		Dependent	
Name	SSN	EmpSSN	DepName
John	999999999	999999999	Emily
Tony	777777777	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

CSE 414 - Spring 2018                      16

### Renaming

- Changes the schema, not the instance

$\rho_{B1, \dots, Bn}(R)$

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

CSE 414 - Spring 2018                      17

### 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**

CSE 414 - Spring 2018                      18

### Natural Join Example

R	A	B
X	Y	
X	Z	
Y	Z	
Z	V	

S	B	C
Z	U	
V	W	
Z	V	

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

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

CSE 414 - Spring 2018 19

### Natural Join Example 2

age	zip	disease
54	98125	heart
20	98120	flu

name	age	zip
Alice	54	98125
Bob	20	98120

$P \bowtie V$

age	zip	disease	name
54	98125	heart	Alice
20	98120	flu	Bob

$\sigma_{\theta} = P.zip = V.zip \text{ AND } P.age = V.age$   
 $\Pi_A = zip, age, \dots$

CSE 414 - Spring 2018 20

### 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  $R \bowtie S$ ?

CSE 414 - Spring 2018 21

### 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:

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

CSE 414 - Spring 2018 22

### Equijoin

- A theta join where  $\theta$  is an equality predicate

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

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

CSE 414 - Spring 2018 23

### Equijoin Example

age	zip	disease
54	98125	heart
20	98120	flu

name	age	zip
p1	54	98125
p2	20	98120

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

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

CSE 414 - Spring 2018 24

## 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  $\theta$
  - No projection
- **Equijoin:**  $R \bowtie_{\theta} S = \sigma_{\theta}(R \times S)$ 
  - Join condition  $\theta$  consists only of equalities
  - No projection
- **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

CSE 414 - Spring 2018

25

## 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

CSE 414 - Spring 2018

26

## 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

CSE 414 - Spring 2018

27

## 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

$P \bowtie J$

$\bowtie_{R \circ J}$   
 $\bowtie_{F \circ J}$

$\bowtie_{L \circ J}$

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

CSE 414 - Spring 2018

28