#### Introduction to Database Systems CSE 414

#### Lecture 10: Relational Algebra

#### Recap: Datalog

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

#### **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: RDMBS internals and query optimization
- Unit 5: Parallel query processing
- Unit 6: DBMS usability, conceptual design
- Unit 7: Transactions

### **Relational Algebra**

Set-based algebra that manipulates relations

– We will extend it to multisets / bags

- In SQL & Datalog we say <u>what</u> we want
- In RA we can express <u>how</u> to get it
- Every DBMS implementations converts a SQL query to RA in order to execute it
- An RA expression is called a <u>query plan</u>

# 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*

#### 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



CSE 414 - Spring 2018

Relational algebra FTW!

### **Relational Algebra Operators**

- Union ∪, intersection ∩, difference -
- Selection σ
- Projection  $\pi$
- Cartesian product X, join M
- (Rename p)
- Duplicate elimination δ
- Grouping and aggregation y
- Sorting τ

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





Only make sense if R1, R2 have the same schema

What do they mean over bags?

### What about Intersection ?

• Derived operator using minus

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

• Derived using join (as we will see later)

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

#### Selection

Returns all tuples which satisfy a condition



- Examples
  - $\sigma_{\text{Salary} > 40000}$  (Employee)
  - $\sigma_{\text{name = "Smith"}}$  (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

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

 Example: project social-security number and names:

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

Different semantics over sets or bags! Why?

## Employee SSN 1234545

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

π <sub>Name,Salary</sub> (Employee)

Name	Salary		Name	Salary
John	20000 ~		John	20000
John	60000		John	60000
John	20000 -	=		

**Bag semantics** 

Set semantics

Which is more efficient?

### **Composing RA Operators**

#### Patient

p4

4

98120

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

98120

no	name	zip	disease			zip	disease
1	p1	98125	flu			98125	flu
2	p2	98125	heart			98125	heart
3	р3	98120	lung			98120	lung
4	p4	98120	heart			98120	heart
σ <sub>disease='heart'</sub> (Patient) т			Π	zip,disease	(o <sub>disease='hea</sub>	<sub>art'</sub> (Patient)	
no	name	zip	disease			zip	disease
2	p2	98125	heart	-	$\rightarrow$	98125	heart

heart

heart

#### **Cartesian Product**

• Each tuple in R1 with each tuple in R2

### R1 × R2

• Rare in practice; mainly used to express joins

#### **Cross-Product Example**

#### Employee

Name	SSN
John	9999999999
Tony	777777777

#### Dependent

EmpSSN	DepName
9999999999	Emily
77777777	Joe

#### **Employee X Dependent**

Name	SSN	EmpSSN	DepName
John	9999999999	9999999999	Emily
John	9999999999	77777777	Joe
Tony	77777777	9999999999	Emily
Tony	77777777	777777777	Joe

CSE 414 - Spring 2018

### Renaming

• Changes the schema, not the instance



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

#### Natural Join



- Meaning:  $R1 \bowtie R2 = \prod_A(\sigma_\theta(R1 \times R2))$
- Where:
  - Selection  $\sigma_{\theta}$  checks equality of all common attributes (i.e., attributes with same names)
  - Projection Π<sub>A</sub> eliminates duplicate common attributes

#### Natural Join Example

S

R

Α	В
Х	Y
Х	Z
Y	Z
Z	V

 B
 C

 Z
 U

 V
 W

 Z
 V

С Α Β **R** ⋈ **S** = Х U Ζ  $\Pi_{ABC}(\sigma_{R.B=S.B}(R \times S))$ Х Ζ V Y Ζ U Y Ζ V Ζ W V

CSE 414 - Spring 2018

#### Natural Join Example 2

#### AnonPatient P

age	zip	disease
54	98125	heart
20	98120	flu

Voters V

name	age	zip
Alice	54	98125
Bob	20	98120

 $\mathsf{P}\bowtie\mathsf{V}$ 

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

#### Natural Join

- Given schemas R(A, B, C, D), S(A, C, E), what is the schema of R ⋈ 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

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

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

 $P \bowtie_{P.zip = V.zip and P.age >= V.age -1 and P.age <= V.age +1} V$ 

### Equijoin

• A theta join where  $\boldsymbol{\theta}$  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?

#### Equijoin Example

#### AnonPatient P

age	zip	disease	
54	98125	heart	
20	98120	flu	

#### Voters V

name	age	zip	
p1	54	98125	
p2	20	98120	

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

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

### Join Summary

- Theta-join:  $R \bowtie_{\theta} S = \sigma_{\theta} (R \times S)$ 
  - Join of R and S with a join condition  $\boldsymbol{\theta}$
  - Cross-product followed by selection  $\theta$
  - No projection
- Equijoin:  $R \bowtie_{\theta} S = \sigma_{\theta} (R \times S)$ 
  - Join condition  $\boldsymbol{\theta}$  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

### 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	

P - X

K oJ R oJ FOJ

#### AnnonJob J

job	age	zip
lawyer	54	98125
cashier	20	98120

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

CSE 414 - Spring 2018