# L21: Relational Algebra

CS3200 Database design (fa18 s2)

https://northeastern-datalab.github.io/cs3200/

Version 11/26/2018

## Our next focus

- The Relational Model
- Relational Algebra
- Relational Algebra Pt. II [Optional: may skip]

## 1. The Relational Model & Relational Algebra

## Algebra

- Algebra is the study of mathematical symbols and the rules for manipulating these symbols;
- e.g., Linear Algebra
- e.g., Relational Algebra
- e.g., Boolean Algebra
- e.g., Abstract algebra (groups, rings, fields, ...)
- e.g., Elementary algebra



## What you will learn about in this section

- The Relational Model
- Relational Algebra: Basic Operators
- Execution

#### Motivation

The Relational model is **precise**, **implementable**, and we can operate on it (query/update, etc.)

Database maps internally into this *procedural language.* 

## A Little History

- Relational model due to Edgar "Ted" Codd, a mathematician at IBM in 1970
  - <u>A Relational Model of Data for Large</u> <u>Shared Data Banks</u>". <u>Communications</u> <u>of the ACM</u> 13 (6): 377–387

- IBM didn't want to use relational model (take money from IMS)
  - Apparently used in the moon landing...



#### Won Turing award 1981



## The Relational Model: Schemata

• Relational Schema:



**Relation name** 

*String, float, int, etc.* are the <u>domains</u> of the attributes

Attributes

### The Relational Model: Data

Stu	de	nt

An attribute (or column) is a typed data entry present in each tuple in the relation

sid	name	gpa			
001	Bob	3.2			
002	Joe	2.8			
003	Mary	3.8			
004	Alice	3.5			
The number of attributes is the <b>arity</b> of					

the relation

### The Relational Model: Data

#### Student

sid	name	gpa
001	Bob	3.2
002	Joe	2.8
003	Mary	3.8
004	Alice	3.5

The number of tuples is the <u>cardinality</u> of the relation

A <u>tuple</u> or <u>row</u> (or *record*) is a single entry in the table having the attributes specified by the schema

#### The Relational Model: Data

#### **Student**

sid	name	gpa
001	Bob	3.2
002	Joe	2.8
003	Mary	3.8
004	Alice	3.5

Recall: In practice DBMSs relax the set requirement, and use <u>multisets</u> (or <u>bags</u>).

A <u>relational instance</u> is a *set* of tuples all conforming to the same *schema* 

#### To Reiterate

• A <u>relational schema</u> describes the data that is contained in a <u>relational</u> <u>instance</u>

Let  $R(f_1:Dom_1,...,f_m:Dom_m)$  be a <u>relational schema</u> then, an <u>instance</u> of R is a subset of  $Dom_1 \times Dom_2 \times ... \times Dom_n$ 

In this way, a <u>relational schema</u> R is a **total function from attribute** names to types

### One More Time

 A <u>relational schema</u> describes the data that is contained in a <u>relational</u> <u>instance</u>

A relation R of arity t is a function: R :  $Dom_1 x \dots x Dom_t \rightarrow \{0,1\}$  *I.e. returns whether or not a tuple of matching types is a member of it* 

Then, the schema is simply the *signature* of the function

Note here that order matters, attribute name doesn't... We'll (mostly) work with the other model (last slide) in which **attribute name matters, order doesn't!** 

## A relational database

- A <u>relational database schema</u> is a set of relational schemata, one for each relation
- A <u>relational database instance</u> is a set of relational instances, one for each relation

Two conventions:

1. We call relational database instances as simply *databases* 

2. We assume all instances are valid, i.e., satisfy the *domain constraints* 

### A Course Management System (CMS)

#### • Relation DB Schema

- Students(sid: string, name: string, gpa: float)
- Courses(cid: string, cname: string, credits: int)
- Enrolled(sid: string, cid: string, grade: string)

Note that the schemas impose effective <u>domain /</u> <u>type constraints</u>, i.e. Gpa can't be "Apple"

Sid	Name	Gpa		Relatio	n		cid	cname	cree
101	Bob	3.2		Instanc	es		564	564-2	4
123	Mary	3.8					308	417	2
	Students		sid	cid	Grad	le		Courses	
			123	564	A				

## 2nd Part of the Model: Querying

SELECT S.name
FROM Students S
WHERE S.gpa > 3.5;

*"Find names of all students with GPA > 3.5"* 

We don't tell the system *how* or *where* to get the data- **just what we want**, i.e., Querying is <u>declarative</u>

To make this happen, we need to translate the *declarative* query into a series of operators... we'll see this next!



Actually, I showed how to do this translation for a much richer language!

## Virtues of the model

- Physical independence (logical too), Declarative
- Simple, elegant clean: Everything is a relation
- Why did it take multiple years?
  - Doubted it could be done efficiently.



## 2. Relational Algebra

### **RDBMS** Architecture

• How does a SQL engine work ?



Declarative query (from user) Translate to relational algebra expression Find logically equivalent- but more efficient- RA expression Execute each operator of the optimized plan!

### **RDBMS** Architecture

• How does a SQL engine work ?



Relational Algebra allows us to translate declarative (SQL) queries into precise and optimizable expressions!

## Relational Algebra (RA)

- Five basic operators:
  - 1. Selection:  $\sigma$
  - 2. Projection:  $\Pi$
  - 3. Cartesian Product: ×
  - 4. Union:  $\cup$
  - 5. Difference: -
- Derived or auxiliary operators:
  - Intersection, complement
  - Joins (natural, equi-join, theta join, semi-join)
  - Renaming: ρ
  - Division

#### We'll look at these first!

And also at one example of a derived operator (natural join) and a special operator (renaming)

## Keep in mind: RA operates on sets!

- RDBMSs use <u>multisets</u>, however in relational algebra formalism we will consider <u>sets</u>!
- Also: we will consider the <u>named perspective</u>, where every attribute must have a <u>unique name</u>
  - $\rightarrow$  attribute order does not matter...

Now on to the basic RA operators...

# 1. Selection ( $\sigma$ )

- Returns all tuples which satisfy a condition
- Notation:  $\sigma_c(R)$
- Examples
  - $\sigma_{\text{Salary} > 40000}$  (Employee)
  - $\sigma_{name = "Smith"}$  (Employee)
- The condition c can be =, <, ≤, >, ≥, <>

Students(sid,sname,gpa)

SQL:





Another example:

SSN	Name	Salary
1234545	John	200000
5423341	Smith	600000
4352342	Fred	500000

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

SSN	Name	Salary
5423341	Smith	600000
4352342	Fred	500000

# 2. Projection $(\Pi)$

- Eliminates columns, then removes duplicates (set perspective!)
- Notation:  $\Pi_{A1,...,An}(R)$
- Example: project social-security number and names:
  - $\Pi_{SSN, Name}$  (Employee)
  - Output schema: Answer(SSN, Name)





Another example:

SSN	Name	Salary
1234545	John	200000
5423341	John	600000
4352342	John	200000

 $\Pi_{SSN}$  (Employee)





Another example:

SSN	Name	Salary
1234545	John	200000
5423341	John	600000
4352342	John	200000

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



Name	Salary
John	200000
John	600000

Note that RA Operators are Compositional!

Students(sid,sname,gpa)

```
SELECT DISTINCT
    sname,
    gpa
FROM Students
WHERE gpa > 3.5;
```

How do we represent this query in RA?

 $\Pi_{sname,gpa}(\sigma_{gpa>3.5}(Students))$ 

## $\sigma_{gpa>3.5}(\Pi_{sname,gpa}(Students))$

Are these logically equivalent?

## 3. Cross-Product (X)

- Each tuple in R1 with each tuple in R2
- Notation:  $R1 \times R2$
- Example:
  - Employee × Dependents
- Rare in practice; mainly used to express joins

Students(sid,sname,gpa)
People(ssn,pname,address)

SQL:

SELECT \*
FROM Students, People;



RA: Students × People



#### Another example: People

ssn	pname	address
1234545	John	216 Rosse
5423341	Bob	217 Rosse

#### Students

sid	sname	gpa
001	John	3.4
002	Bob	1.3

#### Students × People



X

ssn	pname	address	sid	sname	gpa
1234545	John	216 Rosse	001	John	3.4
5423341	Bob	217 Rosse	001	John	3.4
1234545	John	216 Rosse	002	Bob	1.3
5423341	Bob	216 Rosse	002	Bob	1.3

# 4. Renaming $(\rho)$

- Changes the schema, not the instance
- A 'special' operator- neither basic nor derived
- Notation:  $\rho_{B1,...,Bn}$  (R)
- Note: this is <u>shorthand</u> for the proper form (since <u>names, not</u> <u>order</u> matters!):
  - −  $ρ_{A1 \rightarrow B1,...,An \rightarrow Bn}$  (R)

Students(sid,sname,gpa)



We care about this operator *because* we are working in a *named perspective* 





Students

studId	name	gradePtAvg	
001	John	3.4	
002	Bob	1.3	

## 5. Natural Join (⋈)

- Notation:  $R_1 \bowtie R_2$
- Joins R<sub>1</sub> and R<sub>2</sub> on equality of all shared attributes
  - If  $R_1$  has attribute set A, and  $R_2$  has attribute set B, and they share attributes  $A \cap B = C$ , can also be written:  $R_1 \bowtie_C R_2$
- Our first example of a derived RA operator:
  - Meaning:  $R_1 \bowtie R_2 = \prod_{A \bowtie B} (\sigma_{C=D}(\rho_{C \rightarrow D}(R_1) \times R_2))$
  - Where:
    - The rename  $\rho_{C \rightarrow D}$  renames the shared attributes in one of the relations
    - The selection  $\sigma_{\text{C=D}}$  checks equality of the shared attributes
    - The projection  $\Pi_{\text{A}\,\text{U}\,\text{B}}$  eliminates the duplicate common attributes









Another example:



sid	S.name	gpa	ssn	address
001	John	3.4	1234545	216 Rosse
002	Bob	1.3	5423341	216 Rosse

### Natural Join practice



• 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  $R \bowtie S$ ?

Example: Converting SFW Query -> RA



Students(sid,name,gpa)
People(ssn,name,address)

SELECT DISTINCT
gpa,
address
FROM Students S,
People P
WHERE gpa > 3.5 AND
$S_name = P_name;$

 $\Pi_{gpa,address}(\sigma_{gpa>3.5}(S \bowtie P))$ 

How do we represent this query in RA?