Updated 2/4/2022

Topic 1: Data models and query languages Unit 3: Relational Algebra Lecture 06

Wolfgang Gatterbauer

CS7240 Principles of scalable data management (sp22)

https://northeastern-datalab.github.io/cs7240/sp22/

2/4/2022

Pre-class conversations

- Last class recapitulation
 - including a few solutions for examples from last class
- today: let's keep it interactive
 - algebra
 - relational algebra
 - Codd's theorem

Algebra and the connection to logic and queries

- Algebra
- Relational Algebra
 - Operators
 - Independence
 - Power of algebra: optimizations
- Equivalence RA and safe RC (Codd's theorem)

What is "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., Elementary algebra
 - e.g., Abstract algebra (groups, rings, fields, ...)



What is "Algebra"?

- Abstract algebra: studies algebraic structures, which consist of:
 - A domain (i.e. a set of elements)
 - A collection of operators
 - each of arity d; maps a domain of sequences (x₁,...,x_d) to an element y of its codomain (usually that is also the domain)
 - A set of axioms (or identities) that these operators must satisfy.
 - e.g. commutativity: $x \oplus y \equiv y \oplus x$ or $\bigoplus(x,y) \equiv \bigoplus(y,x)$ or $op(x,y) \equiv op(y,x)$
- Examples:
 - Boolean algebra: ({true,false},{∧,∨,¬})
 - − Ring of integers: (ℤ,{+,·})
 - Relational algebra

ring: set equipped with two binary operations with certain properties like distributivity of multiplication over addition

- The definition of an operator allows for composition:
 - e.g. $op_1(op_2(x), op_1(y, op_4(x, z)))$

Based on material by Benny Kimelfeld and Oded Shmueli for 236363 Database Management Systems, Technion, 2018. Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u>



Function composition





Sources: <u>https://www.coursehero.com/sg/college-algebra/composition-of-functions/</u>, <u>https://upload.wikimedia.org/wikipedia/commons/2/21/Function_machine5.svg</u>, <u>https://en.wikibooks.org/wiki/Algebra/Functions</u>, <u>http://www.statisticslectures.com/topics/compositionoffunctions/</u> Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u>



what is the shortest path from s to t?



Wolfgang Gatterbauer. Principles of scalable data management: https://northeastern-datalab.github.io/cs7240/



what is the shortest path from s to t?

Answer: 5 = 3 + 2

min [a + d, a + e, a + f, a + g, ..., c + g] min[3+2, 3+4, 3+7, 3+8, ..., 6+8]

?



what is the shortest path from s to t?

Answer: 5 = 3 + 2

min [a + d, a + e, a + f, a + g, ..., c + g] min[3+2, 3+4, 3+7, 3+8, ..., 6+8]

 $= \min[a, b, c] + \min[d, e, f, g]$ $\min[3,5,6] + \min[2,4,7,8]$

min[x,y]+z = min[(x+z), (y+z)](+ distributes over min)

(Tropical semiring)

• Semiring (\mathbb{R}^{k} , min, +, ∞ , 0)



what is the shortest path from s to t?

Answer: 5 = 3 + 2

Principle of optimality from Dynamic Programming: *irrespective of the initial state and decision, an optimal solution continues optimally from the resulting state*

> min [a + d, a + e, a + f, a + g, ..., c + g] min[3+2, 3+4, 3+7, 3+8, ..., 6+8]

 $= \min[a, b, c] + \min[d, e, f, g]$ $\min[3,5,6] + \min[2,4,7,8]$

min[x,y]+z = min[(x+z), (y+z)](+ distributes over min)



How many paths are there from s to t?



Wolfgang Gatterbauer. Principles of scalable data management: https://northeastern-datalab.github.io/cs7240/



How many paths are there from s to t?

Answer: $12 = 3 \cdot 4$

(Ring of real numbers)

• Semiring $(\mathbb{R},+,\cdot,0,1)$



How many paths are there from s to t?

Answer: $12 = 3 \cdot 4$

count [a·d, a·e, a · f, a · g, ..., c · g] count[1.1, 1.1, 1.1, 1.1, ..., 1.1]12 = count [a, b, c] · count [d, e, f, g] $count[1,1,1] \cdot count[1,1,1]$

+ $[X,Y] \cdot z = +[X \cdot z, Y \cdot z]$ (• distributes over +)

• Semiring $(S, \bigoplus, \bigotimes, 0, 1)$

Semirings generalize this idea



 \bigoplus [a \otimes d, a \otimes e, a \otimes f, a \otimes g, ..., c \otimes g]

 $= \bigoplus$ [a, b, c] $\otimes \bigoplus$ [d, e, f, g]

 $\bigoplus[X,Y] \otimes z = \bigoplus[X \otimes z,Y \otimes z]$ (\$\overline\$ distributes over \$\overline\$)

A... Adjacency matrix, or Arcs

think of dots as "1"s



How many paths of length 2 are there from 7 to 6?



A... Adjacency matrix, or Arcs



How many paths of length 2 are there from 7 to 6?

A... Adjacency matrix, or Arcs

only diagonals and $7 \rightarrow 6$ are shown



How many paths of length 2 are there from 7 to 6?

A... Adjacency matrix, or Arcs

only diagonals and $7 \rightarrow 6$ are shown



Neutral element ∞ instead of D



A... Adjacency matrix, or Arcs

weights) from 7 to 6?

Example graph taken from "Kepner, Gilbert. Graph algorithms in the language of linear algebra, 2011" <u>https://doi.org/10.1137/1.9780898719918</u> Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u> only diagonals and

 $7 \rightarrow 6$ are shown

The Relational Algebra

- In the relational algebra (RA) the elements are relations
 - A relation is a schema together with a finite set of tuples
- RA has 5 primitive operators:
 - Unary: projection, selection
 - Binary: union, difference, Cartesian product
- Each of the 5 is essential or "independent": we cannot define it using the others
 - We will see what exactly this means and how this can be proved
- In practice, we allow many more useful operators that can be defined by the primitive ones (thus also called derived operators)
 - For example, equi-joins via Cartesian product and selection

Company

<u>cid</u>	CName	StockPrice	Country
1	GizmoWorks	25	USA
2	Canon	65	Japan
3	Hitachi	15	Japan

RA vs other Query Languages (QLs)



- There are some subtle (yet important) differences between RA and other QLs. In RA, ...
 - … can tables have duplicate records?
 - ... are missing (NULL) values allowed?
 - ... is there any order among records?
 2
 - ...is the answer dependent on the domain from which values are taken (not just the database at hand)?

RA vs other Query Languages (QLs)



- There are some subtle (yet important) differences between RA and other QLs. In RA, ...
 - … can tables have duplicate records?
 - (RA vs. SQL)
 - ... are missing (NULL) values allowed?
 - (RA vs. SQL)
 - ... is there any order among records?
 - (RA vs. SQL)
 - ...is the answer dependent on the domain from which values are taken (not just the database at hand)?
 - (RA vs. unsafe RC)

Recall: Virtues of the relational model

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



System R is a database system built as a research project at IBM San Jose Research (now IBM Almaden Research Center) in the 1970's. System R introduced the SQL language and also demonstrated that a relational system could provide good transaction processing performance. again in System R and in Eagle, the big project at Santa Teresa. Nevertheless, what kicked off this work was a key paper by Ted Codd – was it published in 1970 in CACM?

Mike Blasgen: Yes.

Irv Traiger: A couple of us from the Systems Department had tried to read it – couldn't make heads nor tails out of it. *[laughter]* At least back then, it seemed like a very badly written paper: some industrial motivation, and then right into the math. *[laughter]*

Bob Yost: I went over there with several other people – I was in the Advanced Systems Development Division – I remember going over there in about 1970 to see this because we were working with the IMS⁸ guys at the time. We couldn't believe it; we thought it's going to take at least ten years before there's going to be anything. And it was ten years. *[laughter]*

Irv Traiger: So we had this 1970 paper; there were a couple of other papers that Ted had written after that; one on a language called DSL/Alpha⁹, which was based on the predicate calculus. Glenn Bacon, who had the Systems Department, used to wonder how Ted could justify that everybody would be able to write this language that was based on mathematical predicate calculus, with universal quantifiers and existential quantifiers and variables and really, really hairy stuff.

RDBMS Architecture

• How does a SQL engine work ?

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



Algebra and the connection to logic and queries

- Algebra
- Relational Algebra
 - Operators
 - Independence
 - Power of algebra: optimizations
- Equivalence RA and safe RC (Codd's theorem)

Relational Algebra (RA) operators

- Five basic operators:
 - 1. Selection: σ
 - 2. Projection: Π
 - 3. Cartesian Product: ×
 - 4. Union: U
 - 5. Difference: –
- Auxiliary operators (sometimes counted as basic):
 - 6. Renaming: ρ ("rho")
- Derived
 - 7. Joins ⋈ (natural, equi-join, theta join, semi-join)
 - 8. Intersection / complement
 - 9. Division

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

Two perspectives: we focus on the <u>named perspective</u>, where every attribute must have a unique name, thus attribute order does not matter (contrast with vectors)

- Extended RA
 - 1. Duplicate elimination δ
 - 2. Grouping and aggregation γ
 - 3. Sorting **τ**

RDBMSs use <u>multisets (bags)</u>, however in RA we will consider <u>sets</u>

Relational Algebra (RA) operators

- Five basic operators:
 - 1. Selection: σ ("sigma")
 - 2. Projection: Π
 - 3. Cartesian Product: ×
 - 4. Union: U
 - 5. Difference: –
- Auxiliary (or special) operator
 - 6. Renaming: ρ ("rho")
- Derived (or implied) operators
 - 7. Joins ⋈ (natural, theta join, equi-join, semi-join)
 - 8. Intersection / complement
 - 9. Division

1. Selection (σ)



- Notation: $\sigma_{c}(R)$
- Examples
 - Employee(<u>ssn</u>,name,salary)
 - $\sigma_{\text{Salary} > 40000}$ (Employee)
 - $\sigma_{name = "Smith"}$ (Employee)
- The condition c can be comparison predicates =, <, ≤, >, ≥, <> combined with AND, OR, NOT

Student(<u>sid</u>,sname,gpa)





1. Selection (σ)

- Returns all tuples which satisfy a condition
- Notation: $\sigma_{c}(R)$
- Examples
 - Employee(<u>ssn</u>, name, salary)
 - $\sigma_{\text{Salary} > 40000}$ (Employee)
 - $\sigma_{name = "Smith"}$ (Employee)
- The condition c can be comparison predicates =, <, ≤, >, ≥, <> combined with AND, OR, NOT

Student(sid,sname,gpa)





 $\sigma_{gpa > 3.5}$ (Student)

1. Selection example



Employee

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

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

?

1. Selection example



Employee

C	SSN	Name	Salary	
$\left\{ \begin{array}{c} \left(\right) \right\} $	1234545	John	20000	
	5423341	Smith	60000	
	4352342	Fred	50000	

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

SSN	Name	Salary
5423341	Smith	60000
4352342	Fred	50000

Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u>

2. Projection (Π)

- Eliminates columns, then removes duplicates (set perspective!)
- Notation: $\Pi_{A1,...,An}(R)$
- Alternative: Π_{-B1,...,Bn}(R)
 "project away" operator (not standard)
- Example: project on social-security number and names:
 - Employee(<u>ssn</u>, name, salary)
 - $\Pi_{SSN, Name}$ (Employee)
 - Output schema: Answer(SSN, Name)

Student(<u>sid</u>,sname,gpa)

SQL:

SELECT DISTINCT sname, gpa FROM Student





2. Projection (Π)

- Eliminates columns, then removes duplicates (set perspective!)
- Notation: $\Pi_{A1,...,An}(R)$
- Alternative: Π_{-B1,...,Bn}(R)
 "project away" operator (not standard)
- Example: project on social-security number and names:
 - Employee(<u>ssn</u>, name, salary)
 - $\Pi_{SSN, Name}$ (Employee)
 - Output schema: Answer(SSN, Name)

Student(<u>sid</u>,sname,gpa)

SQL:

SELECT DISTINCT sname, gpa FROM Student









Employee

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





Wolfgang Gatterbauer. Principles of scalable data management: https://northeastern-datalab.github.io/cs7240/



Employee

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





Name	
Ciara	



Employee

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



Wolfgang Gatterbauer. Principles of scalable data management: https://northeastern-datalab.github.io/cs7240/



Employee

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

$$\Pi_{\text{name, salary}}$$
 (Employee)

Bag semantics

NameSalaryCiara20000Ciara60000



Employee

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





Bag semantics

Name	Salary
Ciara	20000
Ciara	60000
Ciara	20000

Name	Salary
Ciara	20000
Ciara	60000

Composing RA Operators



Patient

no	name	zip	disease
1	p1	98125	flu
2	p2	98125	heart
3	р3	98120	lung
4	p4	98120	heart

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



zip	disease
98125	flu
98125	heart
98120	lung
98120	heart

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

zip	disease
98125	heart
98120	heart

Composing RA Operators

How do we call what we see on this page / the property of these two operators



Patient

no	name	zip	disease
1	p1	98125	flu
2	p2	98125	heart
3	р3	98120	lung
4	p4	98120	heart

π_{zip,disease}(Patient)



zip	disease
98125	flu
98125	heart
98120	lung
98120	heart

σ_{disease='heart'}(Patient)

t) 🗸	
------	--

σ _{disease='heart'} (π _{zip,disease}	(Patient)) 🧹	
--------------------------------	--------------------------	--------------	--

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

zip	disease
98125	heart
98120	heart

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

Composing RA Operators





Patient

no	name	zip	disease
1	p1	98125	flu
2	p2	98125	heart
3	р3	98120	lung
4	p4	98120	heart

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

zip	disease
98125	flu
98125	heart
98120	lung
98120	heart

odisease='heart' (Patient

t) -		
------	--	--

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

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

zip	disease
98125	heart
98120	heart

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

Logical Equivalece of RA Plans





 $\Pi_A(\sigma_{A=5}(R)) \stackrel{?}{\Leftrightarrow} \sigma_{A=5}(\Pi_A(R))$

Do projection & selection <u>commute</u> in this example?



Logical Equivalece of RA Plans





 $\Pi_A(\sigma_{A=5}(R)) \stackrel{?}{\Leftrightarrow} \sigma_{A=5}(\Pi_A(R))$

Do projection & selection <u>commute</u> in this example?

Yes 😊

$\Pi_B(\sigma_{A=5}(R)) \stackrel{?}{\Leftrightarrow} \sigma_{A=5}(\Pi_B(R))$

what about here?



Wolfgang Gatterbauer. Principles of scalable data management: https://northeastern-datalab.github.io/cs7240/

Logical Equivalece of RA Plans





 $\Pi_A(\sigma_{A=5}(R)) \stackrel{?}{\Leftrightarrow} \sigma_{A=5}(\Pi_A(R))$

Do projection & selection <u>commute</u> in this example?

Yes 😊

R(R) $\Pi_B(\sigma_{A=5}(R)) \stackrel{?}{\Leftrightarrow} \sigma_{A=5}(\Pi_B(R))$

what about here?

No 😳