

Informatics 1: Data & Analysis

Lecture 6: Tuple Relational Calculus

Ian Stark

School of Informatics
The University of Edinburgh

Friday 1 February 2013
Semester 2 Week 3



Data Representation

This first course section starts by presenting two common **data** representation models.

- The *entity-relationship (ER)* model
- The *relational* model

Data Manipulation

This is followed by some methods for manipulating data in the relational model and using it to extract information.

- *Relational algebra*
- *The tuple-relational calculus*
- The query language *SQL*

The State We're In

Relational models

- Relations: Tables matching schemas
- Schema: A set of field names and their domains
- Table: A set of tuples of values matching these fields

Relational algebra

A high-level mathematical language of operations on relational tables. Each operation takes one or more tables, and returns another.

selection σ , projection π , renaming ρ , union \cup , difference $-$,
cross-product \times , intersection \cap and different kinds of join \bowtie

Tuple relational calculus (TRC)

A **declarative** mathematical notation for writing **queries**: specifying information to be drawn from the linked tables of a relational model.

Simple Example

All records for students more than 18 years old

$$\{ S \mid S \in \text{Students} \wedge S.\text{age} > 18 \}$$

The set of tuples S such that S is in the table “Students” and has component “age” least 18.

This is like [list comprehension](#) in Haskell

```
[ s | s <- students, age s > 18 ]
```

and similar constructions in other languages.

All are based on “comprehensions” in set theory

Tuple Relational Calculus Basics

Queries in TRC have the general form

$$\{ T \mid P(T) \}$$

where T is a *tuple variable* and $P(T)$ is a logical formula.

Every tuple variable such as T has a *schema*, like rows in a relational table, with fields and their domains. In practice, the details of the schema are usually inferred from the way T appears in $P(T)$.

A **tuple variable** ranges over all possible **tuple values** matching its schema.

The result of the query

$$\{ T \mid P(T) \}$$

is then the set of all possible tuple values for T such that $P(T)$ is true.

Another Example

Names and ages of all students over 20

$$\{ T \mid \exists S . S \in \text{Students} \wedge S.\text{age} > 20 \\ \wedge T.\text{name} = S.\text{name} \wedge T.\text{age} = S.\text{age} \}$$

The set of tuples T such that there is an S in table “Students” with component “age” at least 20 and where S and T have the same values for “name” and “age”.

- Tuple variable S has schema matching the table “Students”.
- Tuple variable T has (only) fields “name” and “age”, with domains to match those of S .
- Even if S has other fields, they do not appear in T or the overall result.

Formula Syntax

Inside TRC expression $\{T \mid P(T)\}$ the logical formula $P(T)$ may be quite long, but is built up from standard logical components.

- Simple assertions: $(T \in \text{Table})$, $(T.\text{age} > 65)$, $(S.\text{name} = T.\text{name})$, ...
- Logical combinations: $(P \vee Q)$, $(P \wedge Q \wedge \neg Q')$, ...
- Quantification:

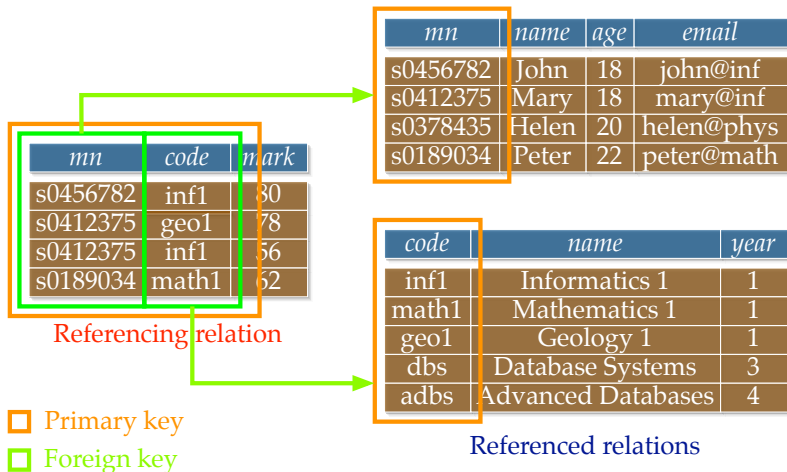
$\exists S . P(S)$ There exists a tuple S such that $P(S)$

$\forall T . Q(T)$ For all tuples T it is true that $Q(T)$

For convenience, we require that for $\exists S . P(S)$ the variable S must actually appear in $P(S)$; and the same for $\forall T . Q(T)$. We also write:

$\exists S \in \text{Table} . P(S)$ to mean $\exists S . S \in \text{Table} \wedge P(S)$

Students and Courses (1/5)



Students and Courses (1/5)

Students taking Informatics 1

$$\{ R \mid \exists S \in \text{Students} . \exists T \in \text{Takes} . \exists C \in \text{Courses} . \\ C.\text{name} = \text{"Informatics 1"} \wedge C.\text{code} = T.\text{code} \\ \wedge T.\text{mn} = S.\text{mn} \wedge S.\text{name} = R.\text{name} \}$$

Schema for S, T and C match those of the tables from which they are drawn. The schema for result R is a single field “name” with string domain, because that’s all that appears here.

One way to compute this in relational algebra:

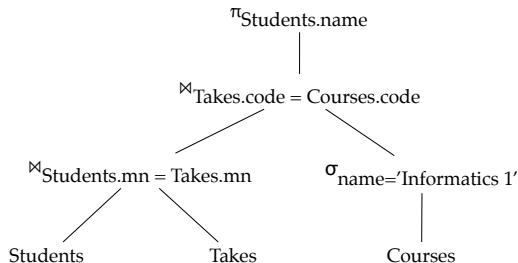
$$\pi_{\text{name}}((\text{Students} \bowtie \text{Takes}) \bowtie (\sigma_{\text{name}=\text{"Informatics 1"}}(\text{Courses})))$$

Relational Algebra

The relational algebra expression can be rearranged without changing its value, but possibly affecting the time and memory needed for computation:

$$\pi_{\text{name}}((\text{Students} \bowtie \text{Takes}) \bowtie (\sigma_{\text{name}=\text{"Informatics 1"}}(\text{Courses})))$$
$$\pi_{\text{name}}(\text{Students} \bowtie (\text{Takes} \bowtie (\sigma_{\text{name}=\text{"Informatics 1"}}(\text{Courses}))))$$
$$\pi_{\text{name}}(\text{Students} \bowtie ((\sigma_{\text{name}=\text{"Informatics 1"}}(\text{Courses})) \bowtie \text{Takes}))$$

We can also visualise this as rearrangements of a tree:



Students and Courses (2/5)

Courses taken by students called "Joe"

$$\{ R \mid \exists S \in \text{Students}, T \in \text{Takes}, C \in \text{Courses} . \\ S.\text{name} = \text{"Joe"} \wedge S.\text{mn} = T.\text{mn} \\ \wedge C.\text{code} = T.\text{code} \wedge C.\text{name} = R.\text{name} \}$$

Note the slightly abbreviated syntax for multiple quantification: we use comma-separated $\exists.., .., ..$ instead of $\exists..\exists..\exists..$

Computing this in relational algebra:

$$\pi_{\text{name}}((\text{Courses} \bowtie \text{Takes}) \bowtie (\sigma_{\text{name}=\text{"Joe"}}(\text{Students})))$$

Students and Courses (3/5)

Students taking Informatics 1 or Geology 1

$$\{ R \mid \exists S \in \text{Students}, T \in \text{Takes}, C \in \text{Courses} . \\ (C.\text{name} = \text{"Informatics 1"} \vee C.\text{name} = \text{"Geology 1"}) \\ \wedge C.\text{code} = T.\text{code} \wedge T.\text{mn} = S.\text{mn} \wedge S.\text{name} = R.\text{name} \}$$

Now the logical formula becomes a little more elaborate.

Computing this in relational algebra:

$$\pi_{\text{name}}((\text{Students} \bowtie \text{Takes}) \bowtie (\sigma_{\text{name}=\text{"Informatics 1"}}(\text{Courses}))) \\ \cup \pi_{\text{name}}((\text{Students} \bowtie \text{Takes}) \bowtie (\sigma_{\text{name}=\text{"Geology 1"}}(\text{Courses}))) \\ \pi_{\text{name}}((\text{Students} \bowtie \text{Takes}) \bowtie (\sigma_{(\text{name}=\text{"Informatics 1"} \vee \text{name}=\text{"Geology 1"})}(\text{Courses})))$$

Students and Courses (4/5)

Students taking both Informatics 1 and Geology 1

$$\{ R \mid \exists S \in \text{Students}, T, T' \in \text{Takes}, C, C' \in \text{Courses} .$$
$$C.\text{name} = \text{"Informatics 1"} \wedge C.\text{code} = T.\text{code} \wedge T.\text{mn} = S.\text{mn}$$
$$C'.\text{name} = \text{"Geology 1"} \wedge C'.\text{code} = T'.\text{code} \wedge T'.\text{mn} = S.\text{mn}$$
$$\wedge S.\text{name} = R.\text{name} \}$$

Computing this in relational algebra:

$$\pi_{\text{name}}((\text{Students} \bowtie \text{Takes}) \bowtie (\sigma_{\text{name}=\text{"Informatics 1"}}(\text{Courses})))$$
$$\cap \pi_{\text{name}}((\text{Students} \bowtie \text{Takes}) \bowtie (\sigma_{\text{name}=\text{"Geology 1"}}(\text{Courses})))$$

Students and Courses (5/5)

Students taking no courses

$$\{ R \mid \exists S \in \text{Students} . S.\text{name} = R.\text{name} \wedge \forall T \in \text{Takes} . T.\text{mn} \neq S.\text{mn} \}$$

Computing this in relational algebra:

$$\pi_{\text{name}}(\text{Students} - \pi_{\text{name,mn}}(\text{Students} \bowtie \text{Takes}))$$

★ Challenge: why not one of these instead?

$$\pi_{\text{name}}(\text{Students} - (\text{Students} \bowtie \text{Takes}))$$

$$\pi_{\text{name}}(\text{Students}) - \pi_{\text{name}}(\text{Students} \bowtie \text{Takes})$$

Relational Algebra vs. Tuple Relational Calculus

Codd gave a proof that relational algebra and TRC are **equally expressive**: anything expressed in one language can also be written in the other.

So why have both?

They give different perspectives and allow the following approach:

- Use relational calculus to specify the information wanted;
- Translate into relational algebra to give a procedure for computing it;
- Rearrange the algebra to make that procedure efficient.

The database language SQL is based on the calculus: well-suited to giving logical specifications, independent of any eventual implementation.

The algebra beneath it is good for rewriting, equations, and calculation.

Domain-Specific Languages



Charles V, 1500–1558
Holy Roman Emperor, King of
Spain, Archduke of Austria

Domain-Specific Languages



Charles V, 1500–1558
Holy Roman Emperor, King of
Spain, Archduke of Austria

“I speak Spanish to God, Italian to women,
French to men and German to my horse.”

Domain-Specific Languages



Charles V, 1500–1558
Holy Roman Emperor, King of
Spain, Archduke of Austria

“I speak Spanish to God, Italian to women,
French to men and German to my horse.”

Attributed, but even Wikipedia is sceptical.

Query Optimization

- ... Rearrange the algebra to make that procedure efficient.

This last part is central to the viability of modern large databases. An effective *query optimizer* will draw up a list of possible *query plans* and compare the costs of all of them, taking account of:

- How much data there is, where it is, how it is arranged;
- What indexes are available, for which tables, and where they are;
- Selectivity: estimates of how many rows a subquery will return;
- Estimated size of any intermediate tables;
- What parts can be done in parallel;
- What I/O and computing resources are available;
- ...