

5 The Relational Data Model: Algebraic operations on tabular data

5.1 Foundation of relational languages

5.2 Relational Algebra operations

5.3 Relational Algebra: Syntax and Semantics

5.4. More Operators

5.5 Special Topics of RA

6 The Relational Data Model: Logic foundation of Data Manipulation

Not presented in class!

Kemper / Eickler: 3.4, 4.6+7; Elmasri / Navathe: chap. 7.4-7.6,
Garcia-Molina, Ullman, Widom: chap. 5, D. Maier Theory of RDB (Online Book -> Lit.)

Context

Part 1: Designing and using database

Database Design:
- developing a relational database schema

Design:
- formal theory

Data handling in relational databases
-Algebra, -Calculus, SQL/DML

Using the Database from application progs

Physical Schema

Part 2: DWH, IR .. Transactions

5.1 Foundation of relational languages



Data Model:

Language for **definition** *and*
handling (manipulation) **of data**

Languages for data handling:

- **Relational Algebra (RA)** as a semantically well defined **applicative language**
- **Relational tuple calculus** (domain calculus):
predicate logic interpretation of data and queries
- **SQL / DML** ('Sequel') – based on **RA and calculus**

SQL: very important in practice

Relational Languages

Goal of DB language design:

Simple and powerful expressions for querying a database

Language should be declarative ("descriptive")

Historically: "Make query formulation 'as easy as in natural language' "

More serious: Queries should be independent of representation of data and implementation aspects (Codd's principle).

Relational Algebra

Algebra objects:

Relations (tables)

$R_1(a_1, \dots, a_n), R_n(b_1, \dots, b_m)$ over domains a_i, b_j, \dots

Algebra operators:

Operators : transform one or more relations into a relation:

$$\tau : R_{i_1}(\dots) \times \dots \times R_{i_k}(\dots) \rightarrow R_{i_m}(\dots)$$

Relational Algebra: only unary and binary operators

Relational Algebra

```
City (name r_id popul .. )
  Oslo    ..  1.31
  Berlin  B   3.47...
  Vienna  ..  957
```

```
Country (c_id r.id capital..)
  'GER'  'B'  Berlin
  'AU'   'V'   Vienna
  ...
```

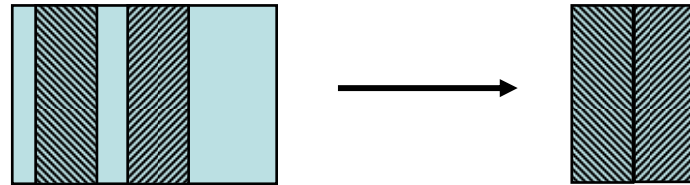
"Name and population of capital of 'Germany' "



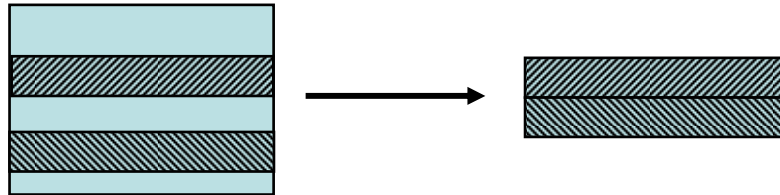
NAME	POPULATION
-----	-----
Berlin	3472009

Basic Operations informally (from chapter 4)

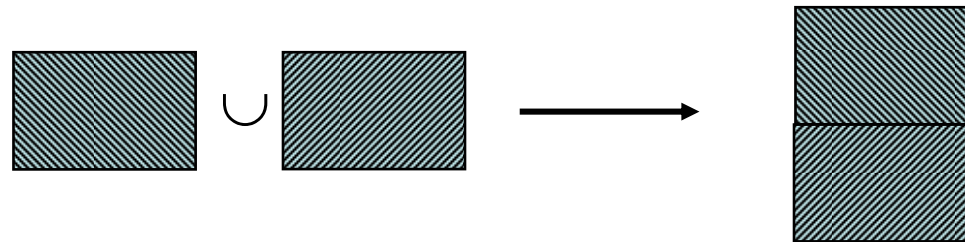
Projection



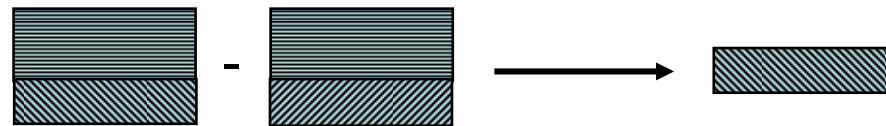
Selection



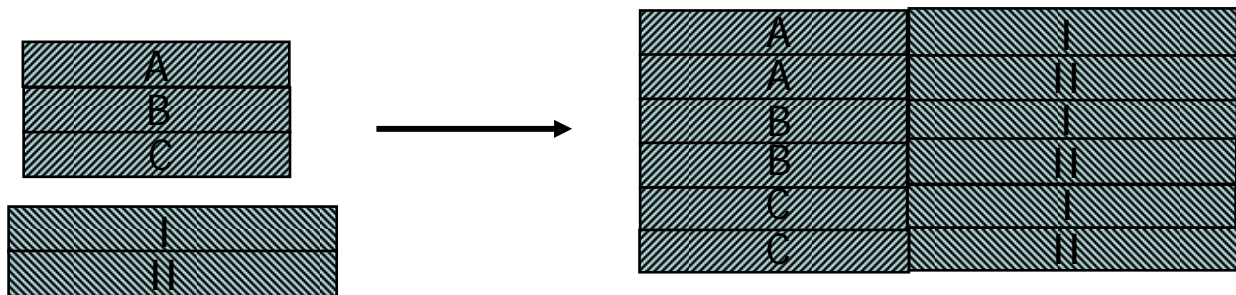
Union



Difference



Cross Product



© HS-2010 (Renaming)

Relational Algebra Operators

{ **Projection** π ,
(Extended) **Cross product** \times ,
Selection σ ,
Union \cup ,
Set difference \setminus ,
Renaming ρ }

is a **base of relational operators**.

Other operators like **join** (\bowtie) can be **expressed**
by $\pi, \sigma, \rho, \times, \cup$

All operators can be expressed in SQL

5.2 Relational Algebra operations

Basic terminology (rep. from above)

Universal set of attributes A , $a_i \in A$ has domain $D(a_i)$

Relation Schema: named n -tuple of attributes

$$RS = R(a_1, \dots, a_n), \{a_1, \dots, a_n\} \subseteq A$$

Schema operator Σ applied to relation R results in the type signature of R : $\Sigma(R) = R_A$

Relational Database Schema: set of relation schemas

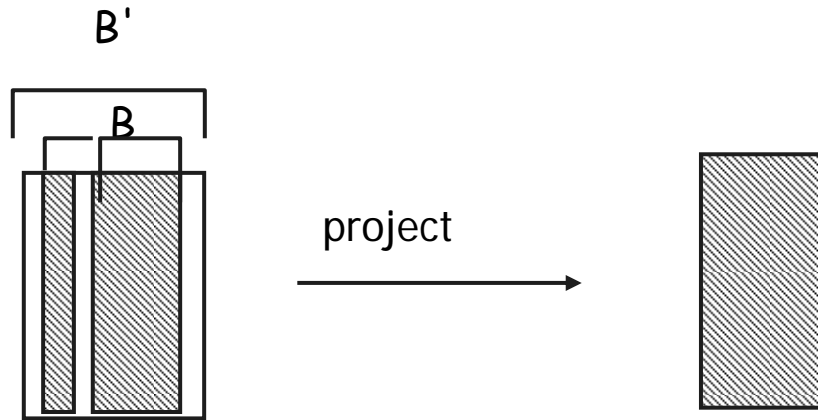
Database Relation R : subset of $D(a_{i_1}) \times \dots \times D(a_{i_n})$

Projection

Let $\Sigma(R) = B'$, $B \subseteq B'$

Projection $\pi_B(R)$ of R on B :

Set of rows from R with the columns not in B eliminated



No duplicates in $\pi_B(R)$ (in theory!)

Def.: $\pi_B(R) = \{r \text{ restricted to } B \mid r \in R\}$
 $= \{r' \mid \text{there is a tuple } r \in R \text{ such that}$
 $r' \text{ is the restriction of } r \text{ to the attributes in } B\}$

Projection (2)

Properties of projection:

- $|R| \geq |\pi_B(R)|$, $B \subseteq \Sigma(R)$
- B contains a key of $R \Rightarrow |\pi_B(R)| = |R|$

Useful for estimating the size of query results
Important for optimization.

SQL equivalent:

```
SELECT DISTINCT  $b_1, b_2, \dots, b_n$  FROM R
```

Extended Cross Product X

Def.: (Extended) Cross product $R \times S =$
 $\{(a_1, \dots, a_n, b_1, \dots, b_m) \mid (a_1, \dots, a_n) \in R, (b_1, \dots, b_m) \in S\}$

R (a1 a2)
1 'A'
5 'Z'

 X

S (b1 b2)
3 'A'
1 'B'

=

T ((a1 , a2) (b1 b2))
(1 'A') (3 'A')
(5 'Z') (3 'A')
(1 'A') (1 'B')
(5 'Z') (1 'B')

SQL equivalent:
SELECT ...
FROM R,S

SELECT ...
FROM
R CROSS JOIN S

Renaming ρ

Def.: Renaming

Attributes: if $\Sigma(R) \cap \Sigma(S) \neq \emptyset$

$\rho \langle \text{attrname} \rangle \leftarrow \langle \text{newAttrname} \rangle \quad (\langle \text{relname} \rangle)$

Relations :

$\rho \langle \text{newname} \rangle (\langle \text{relname} \rangle)$

R	($a1$		$a2$)
		1		'A'	
		5		'Z'	

X

S	($b1$		$a2$)
		3		'A'	
		1		'B'	

$$\Sigma(\rho_{a2 \leftarrow b2}(S)) = \{b1, b2\}$$

Dot notation $R.a2$, $S.a2$ or explicit renaming

Renaming

$\rho_{\langle\text{newname}\rangle}(\langle\text{relname}\rangle)$

Relation $\langle\text{relname}\rangle$ is renamed to $\langle\text{newname}\rangle$ in the context of expression

$\rho_{\langle\text{attrname}\rangle} \leftarrow \langle\text{newAttrname}\rangle(\langle\text{relname}\rangle)$

Attribute $\langle\text{attrname}\rangle$ of relation $\langle\text{relname}\rangle$ is renamed to $\langle\text{newAttrname}\rangle$ in the context of expression

```
 $\pi_{\text{sub.name}}(\sigma_Q(\sigma_P(\text{Employee} \times (\rho_{\text{Sub}}(\text{Employee}))))))$   
where P = "Employee.name = 'Miller' "  
Q = "Sub.boss = Employee.id "
```

```
SELECT ... FROM Employee, Employee Sub  
WHERE ...
```

Set operations

r	(a1	a2)
	1	'B'
	5	'A'
	6	'B'

∪

s	(a1	a2)
	'A'	2
	'B'	5

?

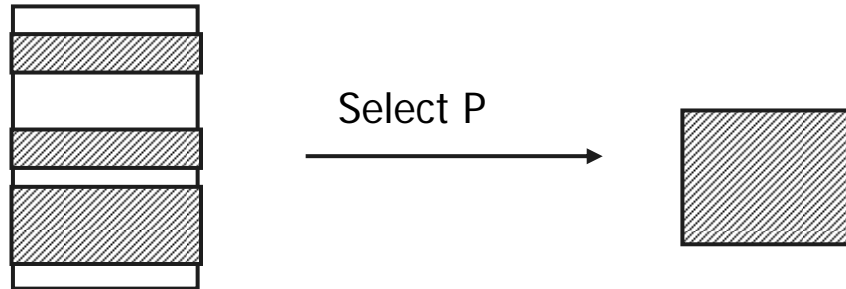
Def.: R and S are called **union-compatible** if the domains of $\Sigma (R) = \Sigma (S)$

R , S union-compatible,
then **set union** and **set difference**
 $R \cup S$ and **$R \setminus S$** as defined on **mathematical sets**

```
SELECT .... FROM R ...  
{UNION | EXCEPT| INTERSECT}  
SELECT ... FROM S ...
```

Selection σ

"Find cities with population more than 1 Mill .



Selection of tuples from a table R according to a predicate P defined on R

Def.: **Selection $\sigma_P(R)$**

Row predicate $P :: R \rightarrow \{\text{true}, \text{false}\}$

$\sigma_P(R) = \{r \mid r \in R, P(r) = \text{true}\}$

Row predicates:

inductively defined by primitive predicates and boolean operators and, or not

Def.: Primitive (simple) predicates

Let a, b be attributes, w value from $\text{dom}(a)$
 $a \theta b$ and $a \theta w$ are primitive predicates
where $\theta \in \{=, \neq, <, \leq, >, \geq\}$

Primitive predicates **compare**
either an attribute and a value
or two attributes

Syntax for **(row) predicates**

- (i) Primitive predicates are predicates
- (ii) If Q , Q' are predicates,
then $Q \wedge Q'$, $Q \vee Q'$ and $\neg(Q)$
are predicates
- (iii) Operator preference and brackets as usual
- (iv) There are no other predicates

"Find countries with more than 5 Mill population
and $GNP \leq 500$ "

Selection of rows

Note:

Selection operator selects the row with all attributes:

$$\Sigma(R) = \Sigma(\sigma_P(R))$$

Size of result depends on **selectivity** of P

$$\text{selectivity} := |\sigma_P(R)| / |R|$$

important for optimization

SQL equivalent (but dupl.):
SELECT ... FROM R
WHERE <row predicate P >

Note:

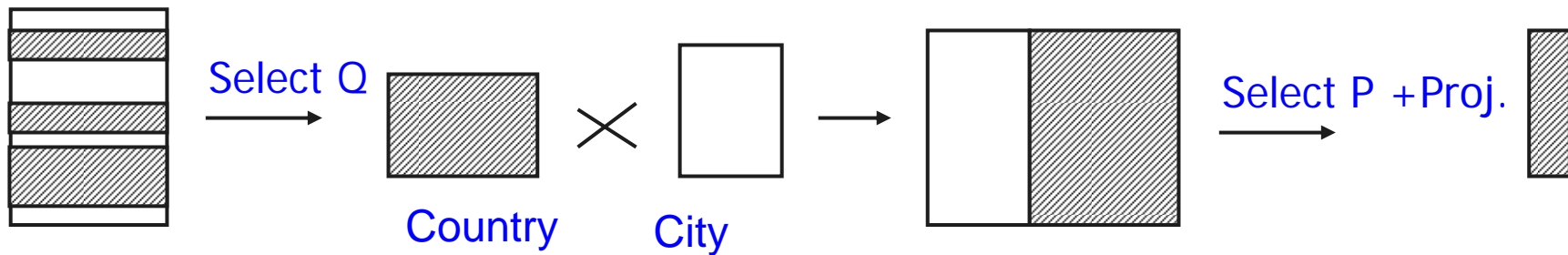
SQL block allows
to combine

π , \times , σ

Combining operators

Country(C_id,name,...,population, ..)

'Find Countries which only consist of its capital and the population > 10000'
(Monaco is an example, Vatican not)



$\pi_{name} (\sigma_P (\sigma_Q (\text{country} \times \text{City})))$
 where Q = "population > 10000 "
 P = "capital = City.name \wedge
 Country.pop = City.pop $\wedge \dots \wedge$

```
SELECT Country.name FROM Country, City
WHERE C_id = City.C_id
and capital = City.name and R_id = City.R_id
and Country.population = City.population
and Country.population > 10000
```

SQL99 syntax

```
SELECT Country.name
FROM Country JOIN City ON C_id = City.C_id
                        and capital = City.name
                        and R_id = City.R_id

WHERE Country.population = City.population
AND Country.population > 10000
```

5.3 Relational Algebra: Syntax and Semantics

Syntax of (simple) Relational Algebra defined inductively :

- (1) Each table identifier is a RA expression
- (2) $\rho_A(B)$, $\rho_{s \leftarrow y}(A)$ are RA expressions where
A,B table identifiers, s, v attribute identifiers
- (3) If E and F are RA expressions then
 $\pi_D(E)$, $\sigma_P(E)$, $E \times F$, $E \cup F$, $E \setminus F$ are RA expressions (if
union-compatible etc.)
where $D \subseteq \Sigma(E)$
- (4) These are all RA expressions

Semantics of Relational Algebra

val is a function which assigns to each relational algebra expression a result table:

$$\text{val ('R')} = R$$

"The value of a relation name is the relation (table)"

$$\text{val ('}\tau (E)\text{')} = \tau (\text{val } (E))$$

where τ is some unary rel. Operation like π

"The value of an unary relational operator applied to an relational algebra expression E is the result of applying the operator to the value of E "

$$\text{val ('E } \omega \text{ F')} = \text{val } (E) \omega \text{ val } (F)$$

where ω is some binary operator like \times

"The value of an unary relational operator applied to a relational algebra expression E is the result of applying the operator to the value of E "

Remarks on RA and SQL

- Rewrite rule

$$\sigma_{Q \wedge P}(R) = \sigma_Q(\sigma_P(R))$$

implicitly used for SQL expression:

SELECT... FROM .. WHERE P (WHERE Q))
does not conform to SQL syntax

- RA results are **sets** (relations),
SQL results are **bags** (duplicates allowed)

To eliminate duplicates write:

```
SELECT DISTINCT ... FROM...  
WHERE ...P AND Q ...
```


Renaming

Renaming , why?

Example: `Employee(id, name, boss, ...)`

Find subordinates of 'Miller'

```

$$\pi_{\text{name}} (\sigma_P (\sigma_Q (\text{Employee}) \times \text{Employee} ))$$
  
where P = "Employee.name = 'Miller' "  
Q = "Employee.boss = Employee.id "
```

RA is a **declarative** language: a **name denotes** the same relation / attribute within one expression

Evaluation example: one table – two roles

Employee		
id	name	boss
001	Abel	NULL
002	Bebel	005
004	Cebel	005
005	Miller	001
006	Debel	001
	



Sub		
id	name	boss
001	Abel	NULL
002	Bebel	005
004	Cebel	005
005	Miller	001
006	Debel	001
	

$\rho_{\text{Sub}}(\text{Employee})$

Renaming

Employee			Sub		
id	name	boss	id	name	boss
001	Abel	NULL	001	Abel	NULL
001	Abel	NULL	002	Bebel	004
...
002	Bebel	005	001	Abel	NULL
002	Bebel	005	002	Bebel	005
...
005	Miller	001	001	Abel	Null
005	Miller	001	002	Bebel	005
005	Miller	001	004	Cebel	005
005	Miller	001	005	Miller	001
005	Miller	001	006	Debel	001
...
006	Debel	001	005	Miller	001
006	Debel	001	006	Debel	001

π_{name}

σ_Q σ_P

5.4 Relational Algebra: More Operators



Some operation sequences occur frequently

⇒ define compound operators

Def.: Join (θ -join)

R, S relations, $R \underset{P}{\bowtie} S$

$= \{(a_1, \dots, a_n, b_1, \dots, b_m) \mid P(a_1, \dots, a_n, b_1, \dots, b_m) \text{ is true}\}$

$= \sigma_P (R \times S)$

where P is a (boolean) predicate composed of primitive predicates of the form

$a \theta b$, $a \in \Sigma(R)$, $b \in \Sigma(S)$, $\theta \in \{=, \neq, <, \leq, >, \geq\}$

(P is the join predicate)

Relational Algebra Join

$$R \bowtie_{R.a < S.c \wedge R.b = S.d} S =$$

1	A	2	1	3	A
2	A	2	1	3	A

The result usually does not have a name

R(a b c)			S(a c d)		
1	A	2	1	3	A
2	A	2	2	2	B
3	C	1	1	2	C

Note: exactly the same as taking the set of all pairs of R and S rows and checking the predicate subsequently

R X S

1	A	2	1	3	A
1	A	2	2	2	B
1	A	2	1	2	C
2	A	2	1	3	A
2	A	2	2	2	B
2	A	2	1	2	C
3	C	1	1	3	A
3	C	1	2	2	B
3	C	1	1	2	C

```
SELECT ...
FROM R JOIN S on (R.a < S.d)
AND (R.b = S.d)
WHERE ...
```

Equijoin: equality comparison

- **Most important type of join:** all primitive predicates in P compare **equality of column values** of two rows at a time : $P \equiv \bigwedge R.x_i = S.y_i$, $\{x_i\} \subseteq \Sigma(R)$, $\{y_i\} \subseteq \Sigma(S)$,
- Implements the "values as pointers" concept of RDB for foreign keys, but is more general.

Example using foreign key: Find Country name title of region having $R_id = 'VAR'$

$$\pi_{name} (\text{Country} \bowtie_{c_id=c_id} \sigma_{R_id='VAR'} (\text{Region}))$$

Relational Algebra: renaming attributes

- Renaming **required**, if **identical column names**
- **No canonical projection** of columns if columns are redundant

R (x, y, z)

1	a	11
5	b	12
6	a	12

S(x', y, z')

7	a	23
6	c	15
9	a	3

$$(R \bowtie S) =$$
$$R.y = S.y$$

R (x, y, z, x', y, z')

1	a	11	7	a	23
1	a	11	9	a	3
6	a	12	7	a	23
6	a	12	9	a	3

Relational Algebra: Natural join

Def.: Natural Join $R \bowtie S$:

equijoin over all **literally identical column names** of R and S and **projection of redundant columns**. Join predicate implicit.



$$R \bowtie S = \pi_{\Sigma(R) \cup \Sigma(S)} (\sigma_P (R \times S))$$

where $P \equiv \bigwedge R.x = S.x, \quad x \in \Sigma(R) \cap \Sigma(S)$

SELECT ... FROM R NATURAL JOIN S

Relational algebra: outer join

Motivation: only tuples of S participate in a join

$R \bowtie S$, which have a "counterpart" in R.

```
Customer(c_no, name, f_name, zip, city)
Phones (phoneNo, c_no)
```

"Print telephon list of customers"

```
 $\pi$  name, phoneNo ( Customer  $\bowtie$  Phones )
```

Customers without phoneNo will not appear

Relational Algebra: outer join

Left outer join $R \bowtie_P S$

Includes $(r, \text{NULL}, \dots, \text{NULL})$ – if there is no join partner for $r \in R$

<table style="border-collapse: collapse; margin: auto;"> <tr><td style="padding: 2px 10px;">a</td><td style="padding: 2px 10px;">b</td><td style="padding: 2px 10px;">c</td></tr> <tr><td style="padding: 2px 10px;">1</td><td style="padding: 2px 10px;">A</td><td style="padding: 2px 10px;">2</td></tr> <tr><td style="padding: 2px 10px;">2</td><td style="padding: 2px 10px;">A</td><td style="padding: 2px 10px;">2</td></tr> <tr><td style="padding: 2px 10px;">3</td><td style="padding: 2px 10px;">C</td><td style="padding: 2px 10px;">1</td></tr> </table>	a	b	c	1	A	2	2	A	2	3	C	1	\bowtie $R.a < S.c \wedge R.b = S.d$	<table style="border-collapse: collapse; margin: auto;"> <tr><td style="padding: 2px 10px;">a</td><td style="padding: 2px 10px;">c</td><td style="padding: 2px 10px;">d</td></tr> <tr><td style="padding: 2px 10px;">1</td><td style="padding: 2px 10px;">3</td><td style="padding: 2px 10px;">A</td></tr> <tr><td style="padding: 2px 10px;">2</td><td style="padding: 2px 10px;">2</td><td style="padding: 2px 10px;">B</td></tr> <tr><td style="padding: 2px 10px;">1</td><td style="padding: 2px 10px;">2</td><td style="padding: 2px 10px;">C</td></tr> </table>	a	c	d	1	3	A	2	2	B	1	2	C	=	<table style="border-collapse: collapse; margin: auto;"> <tr><td style="padding: 2px 10px;">1</td><td style="padding: 2px 10px;">A</td><td style="padding: 2px 10px;">2</td><td style="padding: 2px 10px;">1</td><td style="padding: 2px 10px;">3</td><td style="padding: 2px 10px;">A</td></tr> <tr><td style="padding: 2px 10px;">2</td><td style="padding: 2px 10px;">A</td><td style="padding: 2px 10px;">2</td><td style="padding: 2px 10px;">1</td><td style="padding: 2px 10px;">3</td><td style="padding: 2px 10px;">A</td></tr> <tr><td style="padding: 2px 10px;">3</td><td style="padding: 2px 10px;">C</td><td style="padding: 2px 10px;">1</td><td style="padding: 2px 10px;">-</td><td style="padding: 2px 10px;">-</td><td style="padding: 2px 10px;">-</td></tr> </table>	1	A	2	1	3	A	2	A	2	1	3	A	3	C	1	-	-	-
a	b	c																																												
1	A	2																																												
2	A	2																																												
3	C	1																																												
a	c	d																																												
1	3	A																																												
2	2	B																																												
1	2	C																																												
1	A	2	1	3	A																																									
2	A	2	1	3	A																																									
3	C	1	-	-	-																																									

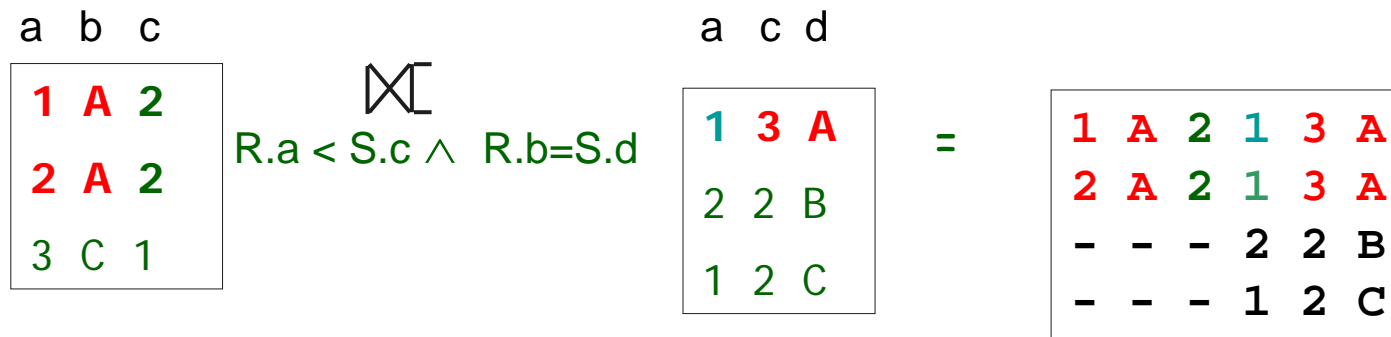
Def.: $R \bowtie_P S =$
 $R \bowtie S \cup \{ (r_1, \dots, r_n, \text{NULL}, \dots, \text{NULL}) \mid (r_1, \dots, r_n) \in R \text{ and}$
 $\text{for all } (s_1, \dots, s_m) \in S: P(r_1, \dots, r_n, s_1, \dots, s_m) = \text{FALSE} \}$

Outer join typically extension of equijoin

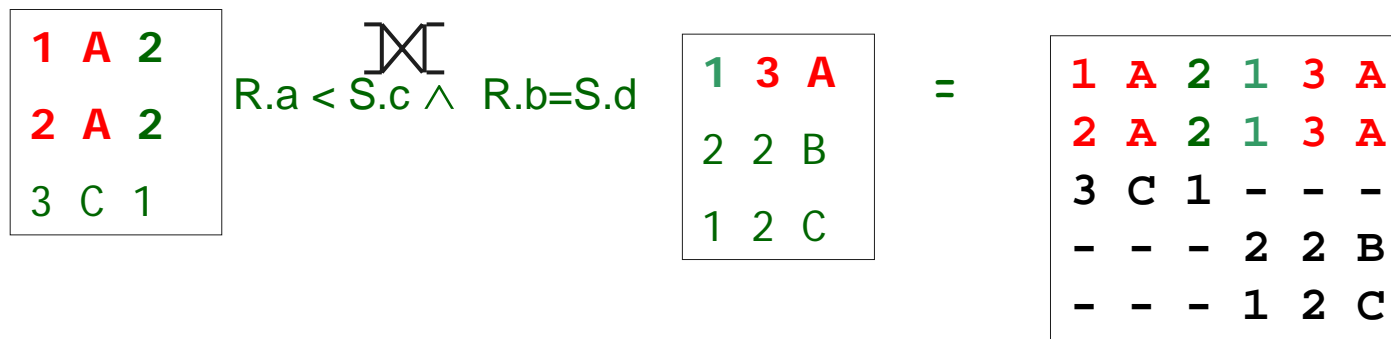
Relational Algebra: outer join

Right outer join $R \bowtie_r S$

Includes (NULL,...NULL, s) – if there is no join partner for $s \in S$



Full outer join: union of left and right outer join





Def.: Semijoin

$$R \bowtie_P S = \Pi_{\Sigma(R)} (R \bowtie S)$$

Left semijoin is the subset of R, each r of which has a corresponding tuple s from S in the join.

Typically extension of equijoin or natural join

R(a b c)

1	A	2
2	A	2
3	C	1

$$\bowtie_{R.a = S.c \wedge R.b = S.d}$$

S(a c d)

1	3	A
2	2	B
1	2	C

=

(a b c)

1	A	2
---	---	---

Right Semijoin defined symmetrically :

$$R \bowtie_S S = \Pi_{\Sigma(S)} (R \bowtie S)$$

Relational Algebra: Base operators

Base

Set of operators which allow to express all other operators

Relational operators

π , σ , \times , \setminus and \cup form a **basis of relational algebra** operators

Means: every RA expression may be expressed only with these operators

Example:
$$R \underset{P}{\bowtie} S = \sigma_P (R \times S)$$

Some rewrite rules for RA

Properties of selection and projection

$$\sigma_P(\sigma_Q(R)) = \sigma_Q(\sigma_P(R))$$

$$\sigma_P(\sigma_P(R)) = \sigma_P(R)$$

$$\sigma_{Q \wedge P}(R) = \sigma_Q(\sigma_P(R)) = \sigma_Q(R) \cap \sigma_P(R)$$

$$\sigma_{Q \vee P}(R) = \sigma_Q(R) \cup \sigma_P(R)$$

$$\sigma_{\neg P}(R) = R \setminus \sigma_P(R), \text{ if } P(r) \text{ defined for all } r \text{ (no NULL!)}$$

$$\text{if } X \subseteq Y \subseteq \Sigma(R) \quad \text{then } \pi_X(\pi_Y(R)) = \pi_X(R)$$

$$\text{if } X, Y \subseteq \Sigma(R) \quad \text{then } \pi_X(\pi_Y(R)) = \pi_{X \cap Y}(R) = \pi_Y(\pi_X(R))$$

$$\text{attr}(P) \subseteq X \subseteq \Sigma(R) \quad \text{then } \pi_X(\sigma_P(R)) = \sigma_P(\pi_X(R))$$

where $\text{attr}(P)$ denotes the set of attributes occurring in P

Relational Algebra operator trees

Algebraic Optimization

- Evaluation of RA expressions in canonical form

$$\pi \dots (\sigma_P (R_1 \times R_2 \times \dots \times R_n))$$

is very inefficient

- How to speed up evaluation of RA (and SQL) expressions?
- Example: Two tables R and S with n and m tuples
Worst case complexity of :

$$\sigma_P (R \bowtie S)$$

is $O(m*n)$

- Interchange of select and join may result in $O(n+m)$

time $\sigma_P (R) \bowtie S$ depending on the join algorithm

Datei Bearbeiten Ansicht Aufgliedern Hilfe

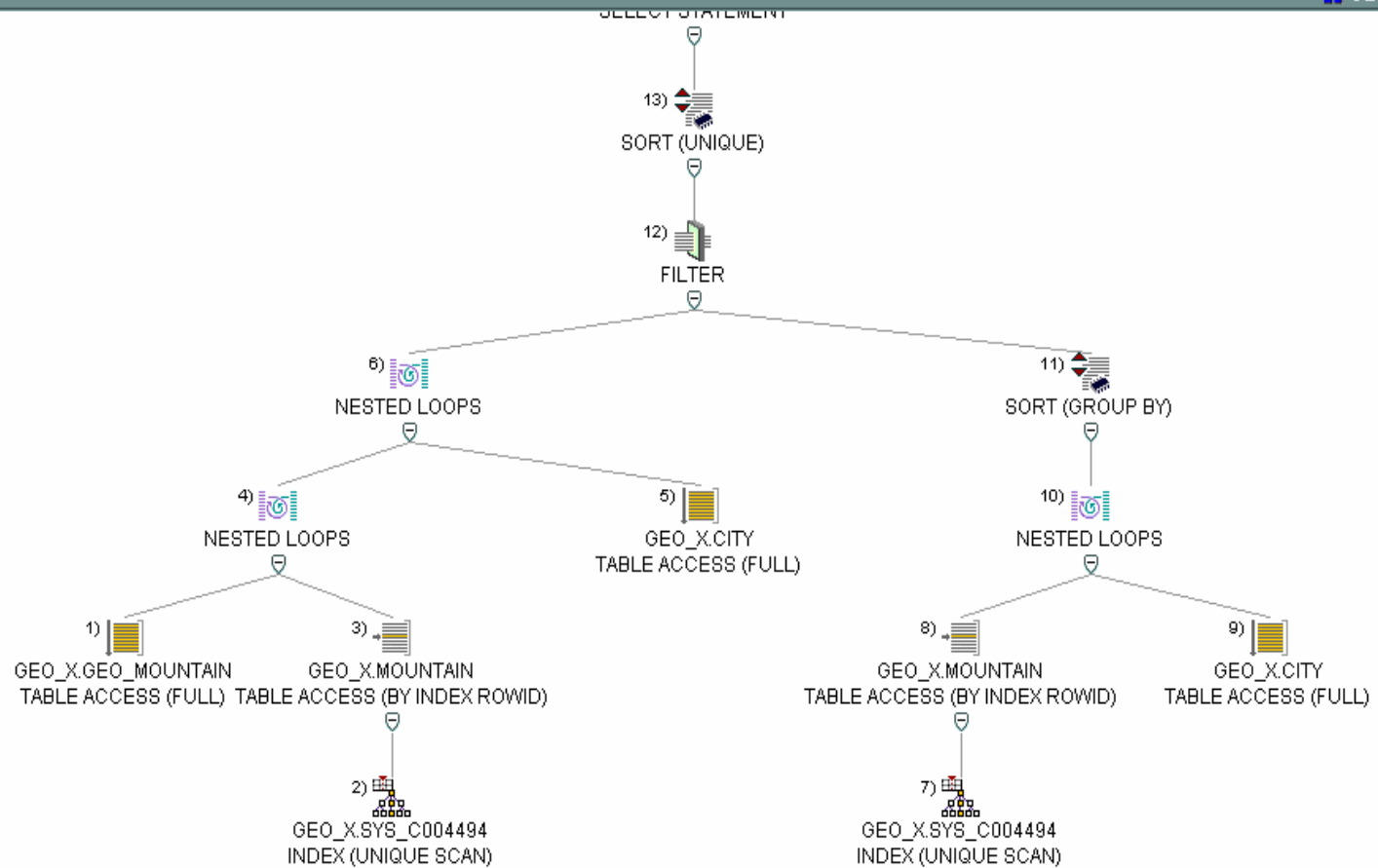
```

SELECT DISTINCT mountain.name, city.name, city.c_id, city.r_id
FROM mountain, city, geoMountain
WHERE geoMountain.c_id = 'SU'
AND mountain.name = geoMountain.mountain
AND (mountain.latitude - city.latitude) * (mountain.latitude -
city.latitude) + (mountain.longit...   ...).longitude) * (
mountain.longitude - city.longitude) = (SELECT min(m.latitude -
city.latitude) * (
m.latitude -

```

Plan:RULE

GEO_X



Dieser Schritt des Plans gibt diese Anweisung als SELECT-Anweisung an.



Row predicates:

P defined over rows (or pairs of rows)

Table predicates

Example: find all countries which are neighbors of all european Countries with population more than 78 Mill

Cannot be answered by comparing individual rows

Predicates with universal quantifier are table predicates

e.g. Find y_0 such that $P(x)$ is true:

$$P(x) \equiv \forall x (\text{PopGT70MillEurope}(x) \Rightarrow (Q(x, y_0)))$$

$$Q(x, y) \equiv x \text{ is neighbor of } y$$

- Express table predicates with base operators?

Relational Algebra: Division

$\text{Course}(\text{id}, \text{title}, \text{semester})$

$T \equiv \text{Course_Stud}(\text{cid}, \text{matr\#})$

ALP4	77
PSem	77
SW	55
SWT	12
SWT	77
ALP4	25
DBS	77
DBS	12

$F \equiv \pi_{\text{id}} (\sigma_{\text{semester}=\text{B4}} \text{Course}(\text{id}, \text{title}, \text{semester}))$

ALP4	B4
DBS	B4
SWT	B4
PSem	B4

Result:

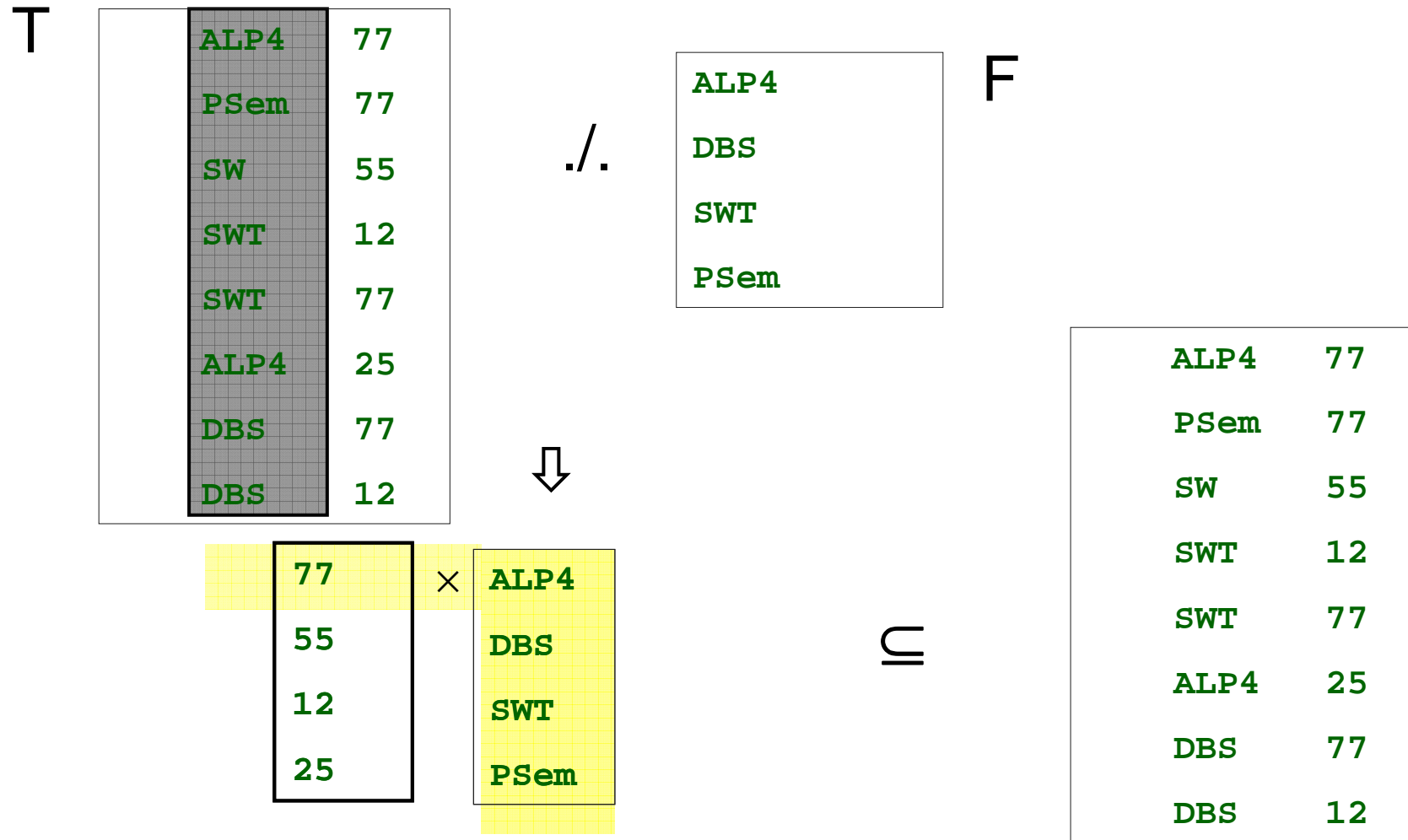
77

Find MatrNo of students who take **all** courses offered for semester B4.

Relational Division

Informally $T \div F$ is the set of all tuples r of T projected on attributes not belonging to F such that $\{(r)\} \times F \subseteq T$

Relational Division: example



$\Rightarrow (77) \in T \div F, \quad \{12, 55, 25\} \notin T \div F$

Def.: Relational Division $T \div F$

Attributes of F are a subset of the attributes of T :

- $\Sigma(F) \subset \Sigma(T)$
- Signature of $T \div F$ is $D = \Sigma(T) \setminus \Sigma(F)$

$$T \div F := \{ t' \mid t' \in \pi_D(T) \wedge (\forall s \in F) (\exists t \in T) \pi_{\Sigma(F)}(\{t\}) = s \wedge \pi_D(\{t\}) = t' \}$$

Simulates a finite "universal" quantification:

"For **all items** x in the table **holds** the predicate **P**"

Relational Algebra Division

$T \div F$ may be defined in terms of other relational operators

$$T \div F = \pi_D(T) \setminus \underbrace{(\pi_D(\pi_D(T) \times F) \setminus T)}_{\text{The "missing" tuples of T}}$$

The "missing" tuples of T

$\underbrace{\hspace{10em}}$
Building the complement

Proof: Assignment

$$D = \Sigma(T) \setminus \Sigma(F)$$

Property of relational division:

Let $D = \Sigma(T) \setminus \Sigma(F)$,

if D contains the key of T and $|F| > 1$ then $T \div F = \emptyset$

5.5. Relational completeness

Completeness

- A DB language L is called **relational complete**, if every RA expression can be expressed in L
- Are there any operations on relations, which cannot be expressed by a finite RA expression (select, project, product or join; **SPJ**) ?
- Yes: **transitive closure** of a relation cannot be expressed in this way

Pred	Descend
Paul	Mary
Mary	Peter
John	Bill
Peter	George

No RA expression to find all decendents of 'Paul'.

Recursion is missing!

What is missing in RA

- Arithmetic operators,
- many practically important operators like **grouping** of results:
"List Students and number of courses they take"

Matr#	NoOfCorses
77	4
55	1
12	2
25	1

ALP4	77
PSem	77
SW	55
SWT	12
SWT	77
ALP4	25
DBS	77
DBS	12

- More Predicates on tables (not rows)

Anyway **relational algebra important conceptual basis** for query languages and query evaluation

RA for optimization

An relational algebra **operator tree** is the data structure representing a RA expression

Algebraic optimization: systematic interchange of operation according to the laws of RA

Does not change time complexity in general, but “makes n small”.

Implementation of Algebraic Optimization by transformation of the operator tree

- Systematic treatment of different optimization techniques → course "DB-Tech"

Summary

Relational algebra: **algebra on tables**

Operators: **project, select, cartesian product, union, set difference, (rename)**

Several compound operators : **join, outer join, semi-join, division**

Serves as a **basis for relational DB languages**

No recursion \Rightarrow not computationally complete

Base of **SQL**

Used for optimization by operator tree transformation

6. The Relational Data Model (*) : Logic foundation of data manipulation - in a nutshell -

- 6.1 Logical foundations of the RDM
- 6.2 Relational Calculus Languages
 - 6.2.1 Tuple calculus
 - 6.2.2 Brief overview of domain calculus
- 6.3 Equivalence of relational languages

Kemper / Eickler: Chap 3.5 , Elmasri/ Navathe: chap. 9.3+9.4
Garcia-Molina, Ullman, Widom: chap. 10,

(*) not discussed in class, not required for exam -

see also reader: logic&databases.pdf

6.1 Logical foundations of the RDM

Predicate logic (PL) view of a DB

Database may be seen as a set of facts:

- $r = (r_1, \dots, r_n) \in R$ for some relation R
- assign a predicate R' to R which is defined:

$$R'(r) = \text{TRUE} \iff r \in R$$

R' is called a **database predicate**

Example:

Movie (25, Amistad, History, 1, Spielberg, 01.05.97)
is a fact, "Movie" is a db predicate

RDM and predicate logic

Restrictions on PL formula

- only database predicates and comparison predicates ($>$, $<$, $=$, \leq , \geq , \neq)
- Variables represent tuples (!)

Open and closed PL formula

- Closed : no free variables, i.e. every variable is bound by a quantifier.
Example: see above
- Open: there are free variables, i.e not closed
- Example:

$\exists t (\text{Tape}(t) \wedge t.\text{movieId} = m.\text{mId} \wedge t.\text{format} = \text{'DVD'})$
Variable m is free in the formula

Open formula as queries

Open formula

$\exists t (\text{Tape}(t) \wedge t.\text{movield} = m.\text{mld} \wedge t.\text{format} = \text{'DVD'})$

An open formula, the free (tuple) variable is m

$\exists m (\text{Movie}(m) \wedge \exists t (\text{Tape}(t) \wedge t.\text{movield} = m.\text{mld} \wedge m.\text{mld} = \text{'4711'} \wedge t.\text{format} = \text{'DVD'}))$

and also

$\exists m (\text{Movie}(m) \wedge \exists t (\text{Tape}(t) \wedge t.\text{movield} = m.\text{mld} \wedge t.\text{format} = \text{'DVD'}))$

are closed and can be evaluated to TRUE | False..

Tuple calculus

Interpret $\{ (s.1) \mid P(x,y,\dots,s) \}$ as:
all rows s which satisfy $P(..)$

$s.1$ means first component of (tuple) variable s

$\{ (s.1) \mid \exists m (\text{Movie}(m) \wedge \exists t (\text{Tape}(t) \wedge t.\text{movieId} = m.\text{mId}$
 $\wedge s.1 = m.\text{title} \wedge t.\text{format} = \text{'DVD'}))$

Formula is open because of $s1$

This is a declarative statement for the set given by the projection of all those s onto the first component, which make the predicate on the right hand side of \mid true.

Open formula as queries

- Implicit requirement: the database predicate of the variables must be known
Technically speaking: the variables must be **range coupled**

Example

- $\{x.3 \mid \text{Movie}(x) \wedge x.\text{title} = P(\dots, x, \dots)\}$
- ' x is a variable, which represents tuples of Movie',
- Query result is the third component of those movie tuples which make P true.

6.2 Calculus Languages

Predicate logic as a query language

Called "calculus" languages for historical reasons

Two types of languages

Domain calculus

All **variables represent typed values** (or domains) of the relations of the DB (domain variables)

Tuple calculus

The **variables** in expressions **represent a row** (tuple) of a relation (tuple variable)

Tuple Calculus language

- Defined over
 - predicates R, S, T, \dots which correspond to database relations
 - set A of attribute names $\{ a, b, \dots \}$
 - values from the domains of A
 - tuple variables r, s, t, \dots

- A tuple calculus expression (with one free variable) has the form

$$\{ \mathbf{s} \mid \mathbf{F}(\mathbf{s}) \}$$

where s is a tuple variable and $F(s)$ a formula in which s occurs free

Example

$\{(s.1, s.2) \mid \exists m (\text{Movie}(m) \wedge m.\text{title} = s.1 \wedge \text{and}$
 $m.\text{year} > '1992' \wedge m.\text{year} = s.2)\}$

- i.e. all movie and year of production titles produced after 1999

Tuple calculus

- s is called the target list
- s is not range coupled and its components are denoted $s.1, s.2, \dots, s.k$, if there are k output components
- $s.i$ has to be connected in F to some range coupled variables

Example:

$\{(s.1, s.2) \mid \exists m (\text{Movie}(m)$
 $\wedge m.\text{title} = s.1 \wedge m.\text{year} > '1992' \wedge m.\text{year} = s.2)\}$

i.e. all movie and year of production titles produced after 1999

Tuple calculus

Simplification

$$\{(s.1, s.2) \mid \exists m (\text{Movie}(m) \wedge m.\text{title} = s.1 \wedge m.\text{year} > '1992' \wedge m.\text{year} = s.2)\}$$

Substitute "logical variables" $s.i$ by row variables

$$\{(m.\text{title}, m.\text{year}) \mid \exists m (\text{Movie}(m) \wedge m.\text{year} > '1992')\}$$

Eliminate existential quantifiers : "all free variables are existentially quantified" ... at least those in target list.

$$\{(m.\text{title}, m.\text{year}) \mid \text{Movie}(m) \wedge m.\text{year} > '1992'\}$$

Range coupling of (row) variable m

Tuple calculus and relational algebra



Tuple calculus expression for algebra operators

- Projection, cross product

$$\pi_{a,b} (R \times S) \equiv \{ (x.a, y.b) \mid R(x) \wedge S(y) \}$$

- Join $R \bowtie_P S$

$$\equiv \{ r.e, \dots, t.k \mid R(r) \wedge S(t) \wedge P \}$$

- Selection

$$\begin{aligned} \text{e.g. } \sigma_{(s.a = v \vee s.a = w) \wedge s.b = s.c} (R) \\ \equiv \{ (s.a, \dots, s.k) \mid R(s) \wedge (s.a = v \vee s.a = w) \wedge s.b = s.c \} \end{aligned}$$

Tuple calculus: examples(1)

Movies (title) **all** copies of which are on loan

$$\{m.title \mid Movie(m) \wedge \forall t (Tape(t) \wedge m.m_id = t.m_Id \\ \Rightarrow \exists x (Rental(x) \wedge x.t_id = t.t_id)) \}$$

Find movie titles available in **all** formats (in the DB)

"... for all formats there exists a tape with this format and this movie"

$$\{m.title \mid Movie(m) \wedge \forall f (Format(f) \Rightarrow \exists x (Tape(x) \wedge \\ f.format = x.format \wedge x.m_id = m.m_id))\}$$

Examples (2)

"Find actors who played together in the same movie."

≡

'There exists an actor and another actor and two different "starring" entries, such that the movie-attributes of both entries are the same and the actor attribute values are the foreign key values for these two actors '

$$\{(a1.stage_name, a2.stage_name) \mid Actor(a1) \wedge Actor(a2) \wedge \exists s1 (Starring(s1) \wedge \exists s2 (Starring(s2) \wedge s1.actor_name = \exists a1.stage_name \wedge s2.actor_name = a2.stage_name \wedge \exists s1.movie_id = s2.movie_id \text{ and } s1 \neq s2) \}$$

Tuple calculus used in Ingres / UC Berkeley as data handling language QUEL. Successor Postgres : SQL

Limitations of RCalc and safe expressions

Limitations, extensions and issues

- Difference to first order predicate logic (FOL)
 - no functions $\forall x (x > 1 \Rightarrow \text{square}(x) > x)$ not allowed
 - FOL interested in formula valid for all domains
e.g. $\forall x P(x) \vee \neg P(x)$
 - RC: Interpretation of tuple calculus expressions over the DB
- What does $\{x / \neg R(x)\} = \{x / \neg \exists t (R(t) \wedge x = t)\}$ mean?
- All tuples NOT belonging to R may not even be a finite set

Safeness of Relational Calculus expressions



**A tuple calculus expression is called safe,
if the result is finite**

- Unfortunately safety property is **not decidable**
- Roughly speaking (syntactically), expressions are safe, if no range variable occurs negated outside an expression which restricts the result set otherwise
- e.g. $\{x \mid R(x)\}$ and $\{x \mid T(x) \wedge \neg R(x)\}$ are safe,
- but $\{x \mid \neg R(x)\}$ is NOT!

6.3 Relational completeness

Relational Algebra and calculus are equivalent

- For each **RA** expression there is an equivalent **safe tuple calculus** expression
- For each **safe tuple calculus** expression there is an **equivalent safe domain calculus** expression
- For each **safe domain calculus** expression there is an **equivalent RA** expression

Equivalent means: results are the same when evaluated over the same DB

- This property of relational languages is called **relational completeness**

Relational complete does not mean computational complete.

Relational completeness

Has been considered as the base line for database query languages: every **query language** should be **as expressive as relational algebra**

SQL is in this tradition, but has introduced many concepts which are difficult or impossible to express in RA

- **grouping and predicates over sets**
e.g. find those movies having the maximal number of copies (DVDs)
- **arithmetic in expressions**, e.g. find cheapest product prices including taxes (in an appropriate DB)
- **partial matches of attribute values**, e.g. find movies the titles of which are LIKE 'To be\$ '
- **application specific comparison functions** (and types),
e.g:
find those customers whose names sound like "Maia"



Relational Algebra

- **Applicative language** on tables for specifying result tables
- **Base for SQL (partially)** and query optimization

Relational Calculus:

- Formal languages syntactically and semantically based on predicate calculus for handling data in relational model
- **Declarative language**, specify *which*, not *how*, data to retrieve
- **Base for QUEL, QBE, SQL (partially)**

I