

Datenbanken
Wintersemester 2019/20
 Prof. Dr. W. May

3. Übungsblatt: SQL

Besprechung voraussichtlich am 8./9.1.2020

Aufgabe 1 (SQL ist relational vollständig) Zeigen Sie, dass SQL *relational vollständig* ist, d.h. zu jedem Ausdruck der relationalen Algebra gibt es einen äquivalenten Ausdruck in SQL.

Das Ziel dieser Aufgabe besteht aus mehreren Ebenen:

- vordergründig Datenbanken: Übersetzung von Ausdrücken der relationalen Algebra nach SQL,
- Prinzipien der Informatik: Strukturelle Induktion, Verwendung von Formalismen.

Zu zeigen ist, dass es für jeden Ausdruck e einen äquivalenten Ausdruck in SQL gibt; d.h., sei \mathbf{R} ein relationales Schema und \mathcal{S} ein Datenbankzustand über \mathbf{R} , dann gibt es eine SQL-Anfrage $q(e)$, so dass das Ergebnis von $q(e)$ (mit den üblichen “Semantik-Klammern” als $\llbracket q(e) \rrbracket_{\mathbf{R}}$ bezeichnet) gleich $\mathcal{S}(e)$ ist.

Um dies für *alle* Algebra-Ausdrücke zu beweisen, folgt man der Definition der Algebra-Ausdrücke durch *strukturelle Induktion*:

Induktionsanfang: die Basisfälle:

- Sei e ein Relationsname r . Dann ist `SELECT * FROM r` die entsprechende SQL-Anfrage $q(e)$.
- Sei $e = \text{cname} : \{x\}$, d.h. die “atomare Tabelle”

cname
x

Dann ist (in Oracle) `SELECT 'x' AS cname FROM DUAL` die gesuchte Anfrage.

Im allgemeinen werden Konstanten aber innerhalb von anderen Anfrage einfach eingesetzt, indem man z.B. `SELECT name, 'Fluss' AS gewaessertyp FROM River` schreibt.

Induktionsschritt: Sei e ein zusammengesetzter Ausdruck, d.h. die Anwendung eines Algebra-operators auf einen oder zwei einfachere Ausdrücke e_1 und ggf. e_2 . Dann existiert nach Induktionsvoraussetzung eine zu e_i äquivalente SQL-Anfrage $q(e_i)$, die man jeweils verwendet.

Fallunterscheidung nach dem jeweiligen Operator (zu betrachten sind nur die Basisoperatoren $\cup, \setminus, \sigma[\dots], \pi[\dots], \rho[\dots]$ und $(\dots) \bowtie (\dots)$ aus der Definition der relationalen Algebra):

Für die ersten 5 Fälle sei das Format von $e_1, q(e_1)$ und ggf. e_2 und $q(e_2)$ jeweils $[A_1, \dots, A_n]$.

- $e = e_1 \cup e_2$: `(SELECT A1, ..., An FROM q(e1))
UNION
(SELECT A1, ..., An FROM q(e2))`
- $e = e_1 \setminus e_2$: `(SELECT A1, ..., An FROM q(e1))
EXCEPT
(SELECT A1, ..., An FROM q(e2))`
- $\pi[A_{i_1}, \dots, A_{i_k}](e_1)$ mit $1 \leq i_j \leq n$: `SELECT Ai1, ..., Aik FROM q(e1)`
Hinweis: in SQL sind außer den A_i auch geeignete Ausdrücke über den Attributnamen erlaubt: Arithmetik, Stringoperationen etc.; z.B.
`SELECT name, population, area, population/area FROM Country.`

- $\sigma[\text{cond}](e_1)$:

```
SELECT A1, ..., An
FROM q(e1)
WHERE cond'
```

(*cond'* ist der *cond* entsprechende Ausdruck in SQL-Syntax)
- $\rho[A_1 \rightarrow C_1, \dots, A_n \rightarrow C_n](e_1)$:

```
SELECT A1 AS C1 , ..., An AS Cn
FROM q(e1)
```
- $e = e_1 \bowtie e_2$, wobei e_2 das Format $[A_{i_1}, \dots, A_{i_k}, B_{k+1}, \dots, B_m]$ hat, $\{i_1, \dots, i_k\} \subseteq \{1 \dots n\}$ und paarweise ungleich:

```
SELECT S.A1, ..., S.An, Bk+1, ... Bm
FROM q(e1) S , q(e2) T;
WHERE S.Ai1 = T.Ai1, ..., S.Aik = T.Aik;
```

Aufgabe 2 (Mondial (SQL)) Gegeben sei folgendes Datenbankschema (Auszug aus Mondial)

```
Country(Name, Code, Capital, Province, Area, Population)
Organization(Name, Abbreviation, Established)
ismember(Organization, Country, Type)
```

Formulieren Sie die folgenden Anfragen in SQL:

(in den Teilaufgaben a) - e) brauchen verschiedene Arten von Mitgliedschaften nicht berücksichtigt werden!)

- Geben Sie von jeder Organisation die Summe der Einwohner aller Mitgliedsländer absteigend geordnet an.
- Welche Länder sind Mitglied in mehr als 60 Organisationen?
- Welche Länder mit einer Fläche von mehr als 500000 km² sind Mitglied in mehr als 60 Organisationen?
- Welche Länder sind in mindestens einer Organisation Mitglied, in der auch Deutschland ('D') Mitglied ist?
- Welche Länder sind in mindestens den Organisationen Mitglied, in denen auch Andorra ('AND') Mitglied ist?
- Zeigen Sie, dass es in der Datenbank keine Organisation gibt, in der alle Länder Mitglied sind!

Diese Anfragen können mit der Web-Schnittstelle zur Mondial-DB getestet werden (siehe Vorlesungsseite).

- Geben Sie von jeder Organisation die Summe der Einwohner aller Mitgliedsländer absteigend geordnet an.

```
SELECT ismember.Organization, SUM(Population)
FROM ismember, Country
WHERE ismember.Country = Country.Code
GROUP BY ismember.Organization
ORDER BY 2 DESC;
```

- Welche Länder sind Mitglied in mehr als 60 Organisationen?

```
SELECT country
FROM ismember
GROUP BY country
HAVING count(*) > 60
```

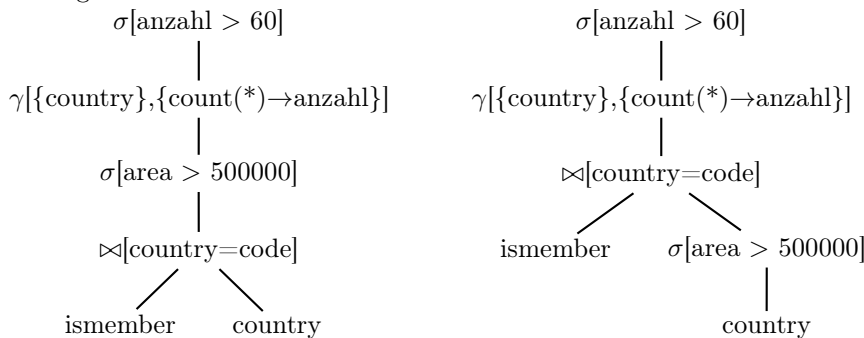
Hinweis: es wird in `ismember` für jedes Land eine Gruppe gebildet. Behalten werden diejenigen Gruppen, die mindestens 60 Einträge (Mitgliedschaften des Landes in einer Organisation) enthalten). Für jede der Gruppen wird eine Zeile im Ergebnis generiert.

c) Welche Länder mit einer Fläche von mehr als 500000 km² sind Mitglied in mehr als 60 Organisationen?

- naheliegendste Formulierung

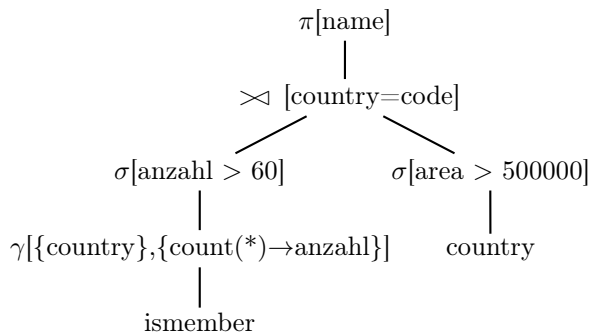
```
SELECT country
FROM ismember i, country c
WHERE i.country = c.code -- join condition
      AND c.area > 500000 -- additional condition
GROUP BY i.country
HAVING count(*) > 60
```

- vgl. Aufgabe zu GROUP-BY-HAVING/Algebra:
als Algebra-Baum:



Auswertungsstrategien:

- ismember (8000 Tupel), an jedes das passende Land (Foreign key → key) dranzoomen (immer noch 8000 Tupel), nach area auswählen (2300 Tupel), gruppieren (50 Gruppen/50 Tupel), Anzahl auswerten (bleiben 6)
- Nur die Länder mit mehr als 500000 km² betrachten (51 Tupel), Mitgliedschaften dranzoomen (2300 Tupel), gruppieren (51 Gruppen), Anzahl auswerten.
- Baum aus Teil b) nehmen: ismember nach country gruppieren, zählen, Anzahl auswerten (18 Tupel). Countries mit mehr als 500000 km² berechnen (50 Tupel). Semijoinen. Bleiben 6 Tupel.



```
SELECT country
FROM (SELECT country
      FROM ismember
      GROUP BY country
      HAVING count(*) > 60),
      (SELECT code FROM country WHERE area > 500000)
WHERE country = code;
```

Letzteres ist sicher keine naheliegende SQL-Formulierung. Es ist aber egal, wie man die Anfrage formuliert, da sie von dem internen Optimierer vorverarbeitet wird.

- See another exercise (Section optimization, indexes)

- d) Welche Länder sind in mindestens einer Organisation Mitglied, in der auch Deutschland ('D') Mitglied ist?

Alternative 1: alles in einem breiten Join:

```
SELECT DISTINCT C.Name
FROM Country C, ismember M1, ismember M2
WHERE C.code = M1.country
      AND M1.organization = M2.organization
      AND M2.country = 'D' AND C.code <> 'D'
```

Alternative 2: "Welches (Land hat eine Mitgliedschaft in einer Organisation), in der auch (Deutschland Mitglied ist)?"

```
SELECT Name
FROM Country, ismember M1
WHERE C1.code = M1.country
      AND M1.organization IN (
          SELECT M2.organization
          FROM ismember M2
          WHERE M2.country = 'D')
      AND C.code <> 'D'
```

Alternative 3: "Für welches Land gibt es eine Mitgliedschaft in einer Organisation, die auch in der Menge der Organisationen, in denen D Mitglied ist, enthalten ist?"

```
SELECT Country.Name
FROM Country
WHERE EXISTS
  (SELECT * from ismember M1
   WHERE M1.country = Country.code
     AND M1.organization IN
       (SELECT M2.organization
        FROM ismember M2
        WHERE M2.country = 'D'))
```

Gegenfrage: welche Länder (die überhaupt irgendwo Mitglied sind) sind nicht in einer Organisation Mitglied, in der Deutschland Mitglied ist?

```
SELECT *
FROM
  (SELECT DISTINCT country AS CC
   FROM ismember)
WHERE NOT EXISTS
  (SELECT *
   FROM ismember M2
   WHERE M2.country='D'
     AND M2.organization IN
       (SELECT organization
        FROM ismember
        WHERE country = CC));
```

- e) Welche Länder sind in mindestens den Organisationen Mitglied, in denen auch Andorra ('AND') Mitglied ist?

-- buggy in Oracle 12

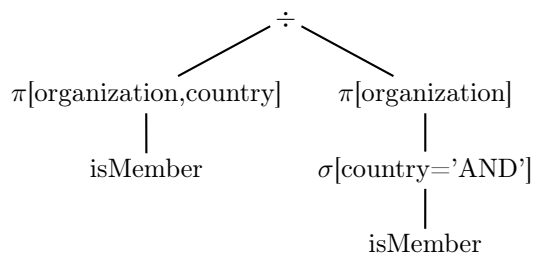
```
SELECT name, code
FROM country c
WHERE NOT EXISTS
  (SELECT * FROM organization o
```

```

WHERE ('AND',o.abbreviation) IN (SELECT country,organization FROM ismember)
AND (c.code,o.abbreviation) NOT IN (SELECT country,organization FROM ismember))
SELECT name, code
FROM country c
WHERE NOT EXISTS
(SELECT organization
FROM ismember i1
WHERE country = 'AND'
AND NOT EXISTS
(SELECT DISTINCT organization
FROM ismember i2
WHERE i1.organization = i2.organization
AND country = c.country ))
SELECT name, code
FROM country c
WHERE NOT EXISTS
((SELECT organization
FROM ismember
WHERE country = 'AND')
MINUS
(SELECT DISTINCT organization
FROM ismember
WHERE country = c.code
))

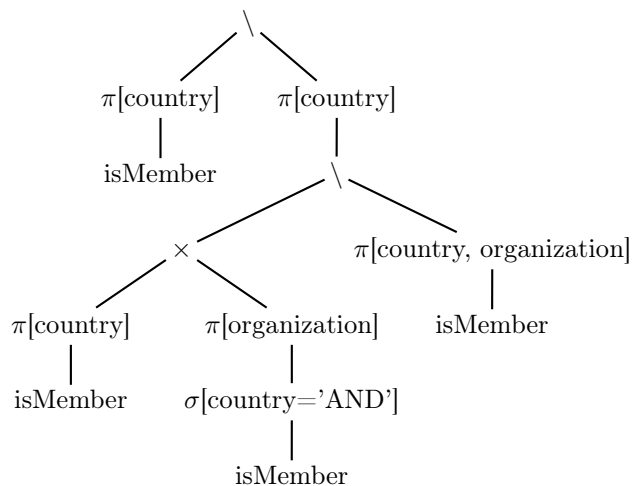
```

Als Algebra-Ausdruck ist es eine relationale Division:



Die Umschreibung ohne \div zeigt die Ähnlichkeit mit dem Pattern $r \div s = \pi[Z](r) - \pi[Z]((\pi[Z](r) \times s) - r)$ für die Division aus der Vorlesung:

Das kartesische Produkt aller Länder mit den Organisationen, in denen Andorra Mitglied ist, bezeichnet die "Soll-Menge", davon werden die tatsächlichen Mitgliedschaften abgezogen, die betreffenden Länder (bei denen Mitgliedschaften "fehlen") durch Projektion ausgewählt, und von der Menge aller Länder abgezogen:



Etwas komplizierter: Mengengleichheit *Welche Länder sind in genau den Organisationen Mitglied, in denen auch Andorra ('AND') Mitglied ist?*

```
SELECT DISTINCT Country
FROM ismember M
  WHERE NOT EXISTS
    ((SELECT DISTINCT organization
      FROM ismember
      WHERE country = 'AND')
     MINUS
     (SELECT DISTINCT organization
      FROM ismember
      WHERE country = M.country))
  AND NOT EXISTS
    ((SELECT DISTINCT organization
      FROM ismember
      WHERE country = M.country)
     MINUS
     (SELECT DISTINCT organization
      FROM ismember
      WHERE country = 'AND'))
```

- f) Zeigen Sie, dass es in der Datenbank keine Organisation gibt, in der alle Länder Mitglied sind!
Alle Organisationen für die es kein Land gibt, das nicht in dieser Organisation ist

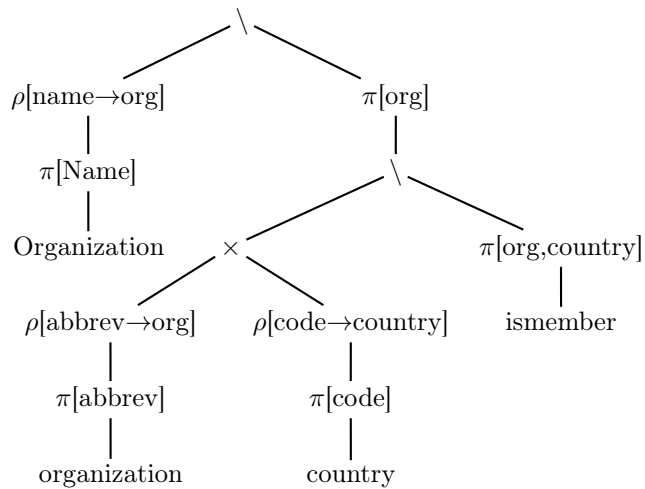
```
SELECT O.Name
FROM Organization O
WHERE NOT EXISTS
  (SELECT name
   FROM country C
   WHERE NOT EXISTS
     (SELECT organization
      FROM ismember I
      WHERE C.code = I.country
      AND I.organization = O.abbreviation))
```

... relationale Division:

$$\pi[\text{org, country}](\text{ismember}) \div \rho[\text{code} \rightarrow \text{country}](\pi[\text{code}](\text{country}))$$

entsprechend der Zerlegung der Division in algebraische Grundoperationen (unter Benutzung

von $\pi[\text{country}](\text{ismember}) = \pi[\text{code}]\text{country}$ und $\pi[\text{org}](\text{ismember}) = \pi[\text{abbrev}](\text{org})$:



Man kann die Anfrage noch dahingehend verfeinern, dass man nur wissen will, ob alle Länder eines Kontinents Mitglied in einer Organisation sind:

```
SELECT abbreviation
FROM organization
WHERE NOT EXISTS
  (SELECT *
   FROM
     (SELECT country CC
      FROM encompasses
      WHERE continent = 'Asia'
     )
   WHERE NOT EXISTS
     (SELECT *
      FROM ismember
      WHERE ismember.organization=organization.abbreviation
        AND ismember.country=CC));
```

Aufgabe 3 (Gruppierung)

- Die Frage nach der größten Landesfläche in der Mondial-Datenbank lautet

```
SELECT MAX(area)
FROM Country;
```

Zusätzlich soll dazu der Landes-Code ausgegeben werden. Warum ist die folgende SQL-Anfrage fehlerhaft? Geben Sie eine entsprechend korrigierte SQL-Anfrage an.

```
SELECT MAX(area), code
FROM Country;
```

- In der Vorlesung wurde für jedes Land die Bevölkerungszahl der größten Stadt ermittelt. Geben Sie eine Anfrage an, die zusätzlich auch den Namen dieser Stadt ausgibt.

- Man könnte sich auf den ersten Blick als Ergebnis der ersten Anfrage `SELECT MAX(area), code FROM Country;` das Tupel (17075200, R) wünschen (Russland ist das größte Land und hat 17075200 km² Fläche). Aber, was sollte dann `SELECT SUM(area), code FROM Country` ergeben? Die Summe der Fläche ist 133287962.24, aber was soll für "code" ausgegeben werden?

Die Aggregationsoperatoren MAX, MIN, SUM, AVG, COUNT nehmen (auf den ersten Blick) als Eingabe eine Menge von Tupeln und wenden den Aggregationsoperator auf die entsprechende Spalte an, und ergeben *einen* atomaren Ergebniswert.

In der obigen falschen Anfrage wird also durch die Anwendung der Aggregatfunktion ein einzeliges Ergebnis für die gesamte Tabelle erzeugt. Das Attribut `code` liefert jedoch nicht ein einzelnes Ergebnis, sondern für jedes Land eines (probieren Sie z.B. mal `SELECT MAX(area), MAX(code) FROM Country;` aus).

Eine korrekte SQL-Anfrage für die Problemstellung ist etwa:

```
SELECT area, code
FROM country
WHERE area =
  (SELECT MAX(area)
   FROM country);
```

- Im obigen Beispiel wird die Aggregatfunktion auf die gesamte aktuelle Relation angewendet. Dies kann durch `GROUP BY attrlist` geändert werden: es werden Mengen von Tupeln gebildet, die in `attrlist` übereinstimmen. Damit werden die Aggregationsoperatoren MAX, MIN, SUM, AVG, COUNT im allgemeinen Fall auf *eine Menge von Mengen von Tupeln* angewendet und werten innerhalb jeder Gruppe den Aggregationsoperator auf der entsprechenden Spalte aus und liefern für jede Gruppe einen atomaren Ergebniswert:

```
SELECT name, country, population
FROM City
WHERE (country, population) IN
  (SELECT country, MAX(population)
   FROM City
   GROUP BY Country);
```

Aufgabe 4 (SQL und Algebra: HAVING) Diese Aufgabe behandelt die `GROUP BY` und `HAVING`-Klauseln von SQL-Anfragen.

- Zeigen Sie: Alle SQL-Anfragen können auch ohne Verwendung von `HAVING` ausgedrückt werden. Geben Sie die SQL-Anfragen mit und ohne `HAVING` für die Anfrage “Welche Länder sind Mitglied in mehr als 60 Organisationen (mit Angabe der Anzahl der Mitgliedschaften)?” an.
- Algebra: Definieren Sie einen Operator `group-by`, der die aus SQL bekannte Funktionalität von `GROUP BY` hat. Betrachten Sie dabei für die Aggregatfunktionen nur einfache Anwendungen auf Attribute, wie z.B. `max(population)`, nicht aber komplexere Ausdrücke wie `max(population/area)`. Gehen Sie dabei wie bei der Definition der Basisoperatoren vor:
 - Welche Parameter müssen dem Operator mitgegeben werden?
 - Welche Signatur hat er?
 - Welche Signatur besitzt die Ergebnisrelation (in Abhängigkeit der Eingaberelation(en))?
 - Wie ist die erhaltene Tupelmenge definiert?
- Geben einen Algebra-Baum für die obige Anfrage “Welche Länder sind Mitglied in mehr als 60 Organisationen?” an.

- Die `HAVING`-Bedingung wird auf die (virtuelle) Relation nach der Gruppierung (und Auswertung der Aggregatoperationen) angewendet. Damit kann man die virtuelle Relation in eine Subquery schachteln:


```

SELECT attrs, op1(attr1), opn(attrn)
FROM tables
WHERE where-cond
GROUP BY group-attrs
HAVING having-cond

SELECT attrs, x1, ..., xn
FROM (SELECT attrs, op1(attr1) as x1, opn(attrn) as xn
      FROM tables
      WHERE where-cond
      GROUP BY group-attrs)
WHERE having-cond mit aliasen x1, ..., xn

```

```

SELECT country, count(*)
FROM ismember
GROUP BY country
HAVING count(*) > 60;

```

```

SELECT country, num
FROM ( SELECT country, count(*) as num
      FROM ismember
      GROUP BY country)
WHERE num > 60;

```

b) **group-by:**

Parameter (Anwendung auf eine Relation $R(\bar{X})$):

- eine Attributmenge $\bar{A} \subseteq \bar{X}$ nach der gruppiert wird,
- eine Menge von Definitionen neuer Attribute. Jede solche Definition enthält den Namen z_i des neuen Attributs sowie dessen Definition. Im angenommenen einfache Fall ist dies ein Paar (Aggregationsoperator, Attribut). Im allgemeinen Fall kann es ein Funktionsausdruck $f_i(\bar{X})$ über \bar{X} sein, wobei f_i Aggregatoperationen (count, sum, max, min, avg) sowie Attribute aus \bar{A} enthalten (z.B. $\max(x_3)$, $\text{sum}(x_4)/\text{count}(*)$, oder $\max(x_1/x_3)$).

Eingabe: eine Relation $R(\bar{X})$.

Ausgabe: eine Relation mit Format $[\bar{A} \cup \bar{Z}]$ wobei $\bar{Z} = \{z_1, \dots, z_n\}$ die Menge der neuen Attributnamen ist.

Definition:

$$\text{group-by}[\bar{A}, (z_1 := f_1(\bar{X}), \dots, z_k := f_k(\bar{X}))](r) = \{t \in \text{Dup}(\bar{AZ}) : t[\bar{A}] \in \pi[\bar{A}](R) \text{ und } t[z_i] = f_i(\sigma[\bar{A} = t[\bar{A}]](r))\}$$

wobei f_i über einer Menge von Tupeln über \bar{X} (der jeweiligen Gruppe) ausgewertet wird.

b) Das Ergebnis der Anwendung von Group-by ist eine ganz “normale” Relation. HAVING entspricht einer Selektion auf dieser.

c) $\pi[\text{country}, \text{num}](\sigma[\text{num} > 60](\text{group-by}[\{\text{country}\}, \{\text{num} := \text{count}(*)\}](\text{ismember})))$

Aufgabe 5 (Relationale Division) Gegeben seien die Relationen $R(A, B)$ und $S(B)$. Prüfen Sie ob der folgende Ausdruck die relationale Division $R \div S$ korrekt abbildet.

```

SELECT A
FROM R
WHERE B IN ( SELECT B FROM S )
GROUP BY A
HAVING COUNT(*) = ( SELECT COUNT (*) FROM S );

```

Geben Sie ggf. die dazu notwendigen Bedingungen oder Korrekturen am Ausdruck an und diskutieren Sie seine Effizienz gegenüber den in der Vorlesung präsentierten Ausdrücken.

Der Ausdruck ist prinzipiell nicht schlecht. Er ist aber nicht korrekt, wenn Duplikate beteiligt sind (echte Duplikate in Tabellen sind selten, aber oft entsteht S oder T durch eine Subquery).

Die folgenden minimalen Gegenbeispiele zeigen das Problem:

- Wenn $R(A, B)$ Duplikate enthält, kann ein A -Wert auch im Ergebnis enthalten sein, wenn er nicht mit allen geforderten B -Werten in R auftritt:

R	
A	B
a	1
a	1

S
B
1
2

- Wenn $R(A, B)$ Duplikate enthält, kann ein A -Wert im Ergebnis fehlen, auch wenn er mit allen geforderten B -Werten in R auftritt:

R	
A	B
a	1
a	1
a	2

S
B
1
2

- Wenn $S(B)$ Duplikate enthält, fehlen Ergebnisse:

R	
A	B
a	1

S
B
1
1

Die Definition der Semantik der Division ist rein mengenorientiert, so dass a hier im Ergebnis erscheinen müsste.

Die obigen Minimalbeispiele erscheinen etwas absurd, kommen aber häufig vor, wenn A oder B durch eine Subquery berechnet sind, z.B.

R'		
A	B	C
a	1	x
a	1	y

oder

S'	
B	D
1	i
1	j

- Korrektur der Anfrage:

```
SELECT A
FROM R
WHERE B IN ( SELECT B FROM S )
GROUP BY A
HAVING COUNT(DISTINCT *) = ( SELECT COUNT (DISTINCT *) FROM S );
```

Für das Beispiel aus der Vorlesung, "Organisationen, die auf allen Kontinenten ein Mitglied haben", ergibt sich

```
SELECT organization
FROM (SELECT organization, continent
      FROM ismember i, encompasses e
      WHERE i.country = e.country)
GROUP BY organization
HAVING COUNT(DISTINCT continent) = (SELECT COUNT (*) FROM continent );
```

- Effizienz: Intern werden üblicherweise (ad-hoc)-Indexe verwendet (in diesem Fall über $\pi[A](R)$ zur Unterstützung des GROUP BY und des COUNT und $\pi[B](S)$). Damit kann diese Variante schneller sein als die normalerweise angegebene, wenn S nicht durch eine korrelierte Unterabfrage berechnet werden muss.
-
-

Ausserdem:

- Formulieren Sie die Anfragen vom vorigen Blatt auch in SQL, und formulieren Sie die Anfragen an Mondial von diesem Blatt soweit möglich auch in der relationalen Algebra.
- Weitere Aufgaben finden Sie auf dem ersten Übungsblatt des SQL-Praktikums (<http://dbis.informatik.uni-goettingen.de/Teaching/DBP-SS19/>)
Dort finden Sie auch detaillierte Folien sowie ein Skript zu SQL ...