

Datenbanken
Wintersemester 08/09
 Prof. Dr. W. May

3. Übungsblatt: SQL

Besprechung am 3.12./10.12.

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.

Gegeben seien die Relationen $R(A_1, \dots, A_n)$, $S(A_1, \dots, A_n)$ und $T(A_{i_1}, \dots, A_{i_k}, B_{k+1}, \dots, B_m)$, $\{i_1, \dots, i_k\} \subseteq \{1 \dots n\}$ und paarweise ungleich.

$R \cup S$:	<pre>SELECT A₁, ..., A_n FROM R UNION SELECT A₁, ..., A_n FROM S;</pre>
$R \setminus S$:	<pre>SELECT A₁, ..., A_n FROM R EXCEPT SELECT A₁, ..., A_n FROM S;</pre>
$\sigma[F](R)$:	<pre>SELECT A₁, ..., A_n FROM R WHERE \bar{F}; \bar{F} ist der F entsprechende Ausdruck in SQL-Syntax</pre>
$\pi[A_1, \dots, A_i](R)$:	<pre>SELECT A₁, ..., A_i FROM R;</pre>
$\rho[A_1 \rightarrow C_1, \dots, A_n \rightarrow C_n](R)$:	<pre>SELECT A₁ AS C₁, ..., A_n AS C_n FROM R</pre>
$S \bowtie T$:	<pre>SELECT A₁, ..., A_n, B_{k+1}, ..., B_m FROM S, T; WHERE S.A_{i₁} = T.A₁, ..., S.A_{i_k} = T.A_k;</pre>

Aufgabe 2 (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 koennt sich auf den ersten Blick als Ergebnis der ersten Anfrage `SELECT MAX(area), code FROM Country;` das Tupel (17075200, R) wünschen (Russland ist das groesste Land und hat 17075200 km² Fläche). Aber, was sollte dann `SELECT SUM(area), code FROM Country` ergeben? Die Summe der Flaeche 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 3 (Zusätzliche Algebraoperatoren) Man kann weitere Algebra-operatoren definieren, die zusätzliche Funktionalität bereitstellen. Definieren Sie zwei Operatoren `group-by` und `top-k` wie folgt:

- `top-k`: Für eine gegebene Relation `country(...,area,...)` enthält `topk(country, area, desc, 5)` die flächenmä"sig größten 5 Länder.
(die Bedeutung des `top-k`-Operators besteht in der Praxis darin, sich von einer komplexen Anfrage die k besten Ergebnisse möglichst schnell zurückzugeben, wobei nach Möglichkeit garnicht alle Ergebnisse wirklich berechnert werden müssen).
- `group-by` soll die aus SQL bekannte Funktionalität von `GROUP BY` haben,

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?

topk: Der top-k-Operator ist eng mit der Selektion verwandt. Allerdings ist das Selektionskriterium keine Bedingung an ein einzelnes Tupel, sondern eine Bedingung, die alle Tupel in R berücksichtigt: Parameter (Anwendung auf eine Relation \bar{X}):

- eine Attributmengende $\bar{A} \subseteq \bar{X}$,
- ein Vergleichsoperator, mit dem Tupel über \bar{A} verglichen/geordnet werden können. Falls \bar{A} einelementig ist, z.B. $\{\text{Area}\}$, ist die Semantik offensichtlich. Wenn \bar{A} mehrere Attribute umfaßt, muß ein Operator "irgendwie" definiert/angegeben werden – dies ist ähnlich zu den MAP/ORDER FUNCTIONS zu Objekttypen in PL/SQL.
Beispiel: Sporttabelle nach Punkten/Tordifferenz/geschossene Tore ordnen.
- k : wieviele Ergebnisse sollen ausgegeben werden?

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

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

Definition:

$$\begin{aligned} \text{topk}[k, \bar{A}, <](r) &= \{t \in r : |\{t' \in R : \pi[\bar{A}](t') > \pi[\bar{A}](t)\}| < k\} \\ &= \{t \in r : |\{\sigma[\bar{A} > \bar{A}(t)](r)\}| < k\} \end{aligned}$$

Hinweis: Wenn mehrere Tupel gleichgroße Werte von \bar{A} haben, können auch $\ell > k$ Tupel ausgegeben werden; für jedes dieser Tupel gibt es aber maximal $k - 1$ größere.

group-by: Parameter (Anwendung auf eine Relation \bar{X}):

- eine Attributmengende $\bar{A} \subseteq \bar{X}$,
- für jede zu aggregierende Spalte y_i von $\bar{X} \setminus \bar{A}$ kann eine Aggregationsoperation op_i sowie ein neuer Name z_i für das Ergebnis angegeben werden.

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

Ausgabe: eine Relation $R((\bar{X} \setminus \bar{A}) \cup \bar{Z})$ mit $\bar{Z} = \{z_1, \dots, z_n\}$.

Definition:

$$\begin{aligned} \text{group-by}[\bar{A}, ((y_1, op_1, z_1), \dots, (y_i, op_i, z_i))](r) \\ = \{t \in \text{Tup}(\bar{A}\bar{Z}) : t[\bar{A}] \in \pi[\bar{A}](R) \text{ und } t[z_i] = op_i(\pi^*[y_i](\sigma[\bar{A} = t[\bar{A}]](r)))\} \end{aligned}$$

wobei π^* keine Duplikateliminierung vornimmt.

Aufgabe 4 (Mengenoperationen) Seien $R(\bar{X}), S(\bar{X})$ Schemata und $\bar{Y} \subseteq \bar{X}$. Zeigen Sie oder widerlegen Sie:

- $\pi[\bar{Y}](R \cup S) = \pi[\bar{Y}](R) \cup \pi[\bar{Y}](S)$
- $\pi[\bar{Y}](R \cap S) = \pi[\bar{Y}](R) \cap \pi[\bar{Y}](S)$
- (falls Anfrageoptimierung und interne Auswertung in der Vorlesung bereits besprochen wurden): wie werden die jeweiligen Ausdrücke am effizientesten ausgewertet?

a) Es gilt Gleichheit:

$$\begin{aligned} \pi[\bar{Y}](R \cup S) &= \{\mu \in \text{Tup}(\bar{Y}) \mid \text{es gibt } \nu \in R \cup S \text{ so dass } \mu = \pi[\bar{Y}](\nu)\} \\ &= \{\mu \in \text{Tup}(\bar{Y}) \mid \text{es gibt } \nu_1 \in R \text{ so dass } \mu = \pi[\bar{Y}](\nu_1)\} \\ &\quad \text{oder es gibt } \nu_2 \in R \cup S \text{ so dass } \mu = \pi[\bar{Y}](\nu_2)\} \\ &= \{\mu \in \text{Tup}(\bar{Y}) \mid \mu \in \pi[\bar{Y}](R) \text{ oder } \mu \in \pi[\bar{Y}](S)\} \\ &= \pi[\bar{Y}](R) \cup \pi[\bar{Y}](S) \end{aligned}$$

b) Diese Gleichung gilt nicht: Betrachte

R	
A	B
a	1
a	2
b	4
c	5
d	7

S	
A	B
a	2
a	3
b	4
c	6
e	8

R ∩ S	
A	B
a	2
b	4

π[A](R)
A
a
b
c
d

π[A](S)
A
a
b
c
e

π[A](R) ∩ π[A](S)
A
a
b
c

c ist im Ergebnis enthalten, obwohl es in R und in S mit unterschiedlichen Werten für B auftritt.
 Sei $\bar{Z} := \bar{X} \setminus \bar{Y}$.

$$\begin{aligned} \pi[\bar{Y}](R \cap S) &= \{\mu \in \text{Tup}(\bar{Y}) \mid \text{es gibt } \nu \in \text{Tup}(\bar{Z}) \text{ so dass } \mu\nu \in R \cap S\} \\ &= \{\mu \in \text{Tup}(\bar{Y}) \mid \text{es gibt } \nu \in \text{Tup}(\bar{Z}) \text{ so dass } \mu\nu \in R \text{ und } \mu\nu \in S\}. \\ \pi[\bar{Y}](R) \cap \pi[\bar{Y}](S) &= \{\mu \in \text{Tup}(\bar{Y}) \mid \mu \in \pi[\bar{Y}](R) \text{ und } \mu \in \pi[\bar{Y}](S)\} \\ &= \{\mu \in \text{Tup}(\bar{Y}) \mid \text{es gibt } \nu_1 \in \text{Tup}(\bar{Z}) \text{ so dass } \mu\nu_1 \in R \text{ und} \\ &\quad \text{es gibt } \nu_2 \in \text{Tup}(\bar{Z}) \text{ so dass } \mu\nu_2 \in S\} \end{aligned}$$

Hierbei kann $\nu_1 \neq \nu_2$ sein.

c) Bei Mengenoperationen (Duplikateliminierung, bei Schnittmenge und Mengendifferenz auch Suche) bietet es sich an, einen Index (Baum oder Hash) zur Laufzeit anzulegen.

a) über R iterieren, jedes Tupel projizieren und Ergebnis in Index einfügen. Dann über S iterieren, jedes Tupel projizieren und Ergebnis einfügen. Alles zusammen auslesen (Hashtabelle: ungeordnet). Falls bereits ein Index über $R[\bar{Y}]$ oder $S[\bar{Y}]$ in der Datenbank existiert, kann man statt Iteration über die Tupel direkt über die Indexteinträge iterieren.

$$|R| \cdot (1 + \text{indextiefe}) + |S| \cdot (1 + \text{indextiefe}) + |\text{Ergebnis}|$$

b) über R iterieren, jedes Tupel in Index einfügen. über S iterieren, jedes Tupel im Index suchen, gefundene Tupel projizieren und das Ergebnis in einen Ergebnis-Index einfügen. Zuletzt diesen ausgeben.

$$|R|y + |S| \cdot (1 + \text{indextiefe}) + |R \cap S| \cdot (1 + \text{ergebnisindextiefe}) + |\text{Ergebnis}|$$

Aufgabe 5 (Mondial (SQL)) Gegeben sei folgendes Datenbankschema (Auszug aus Mondial)

```
Country(Name, Code, Capital, Province, Area, Population)
Organization(Name, Abbreviation, Established)
Is_member(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?

- d) Welche Länder sind in mindestens einer Organisation Mitglied, in der auch Deutschland ('D') Mitglied ist?
- e) Welche Länder sind in mindestens den Organisationen Mitglied, in denen auch Andorra ('AND') Mitglied ist?
- f) 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).

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

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

- b) Welche Länder sind Mitglied in mehr als 60 Organisationen?

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

Hinweis: es wird in `is_member` 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? 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, Is_Member M1, Is_Member 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, Is_Member M1
WHERE C1.code = M1.country
      AND M1.organization IN (
          SELECT M2.organization
          FROM Is_Member 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 is_member M1
```

```

WHERE M1.country = Country.code
AND M1.organization IN
  (SELECT M2.organization
   FROM Is_Member 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 is_member)
WHERE NOT EXISTS
  (SELECT *
   FROM is_member M2
   WHERE M2.country='D'
   AND M2.organization IN
     (SELECT organization
      FROM is_member
      WHERE country = CC));

```

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

```

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

```

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

```

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

```

WHERE country = 'AND'))

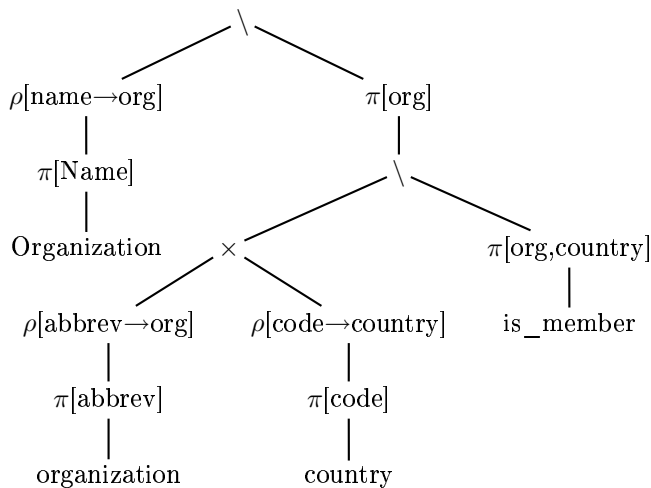
f) Zeigen Sie, dass es in der Datenbank keine Organisation gibt, in der alle Länder Mitglied sind!
Alle Organisationen für die es kein Land gibt, das nicht in dieser Organisation ist

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

... relationale Division:

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

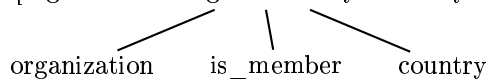
entsprechend der Zerlegung der Division in algebraische Grundoperationen (unter Benutzung von $\pi[\text{country}](\text{is_member}) = \pi[\text{code}]\text{country}$ und $\pi[\text{org}](\text{is_member}) = \pi[\text{abbrev}](\text{org})$):



Man kann auch das NOT EXISTS entsprechend der Vorlesung in die Algebra übersetzen:

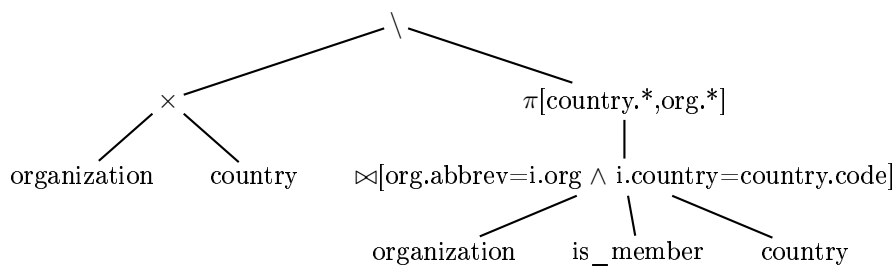
- Innere Subquery: importiert werden *organization* und *country*; also Join mit nachfolgendem von diesen beidem mit *is_member*:

$\bowtie[\text{org.abbrev} = \text{i.org} \wedge \text{i.country} = \text{country.code}]$



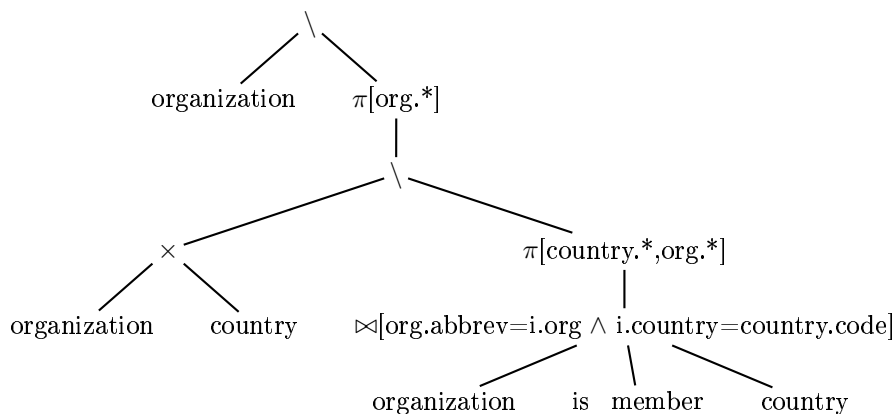
Das Ergebnis dieses Ausdruckes wird auf die Schlüsselattribute der korrelierenden Relationen *organization* und *country* projiziert (ergibt diejenigen Paare von Organisationen und Länder, für die eine Mitgliedschaft existiert).

- Die Grundmenge der mittleren Subquery ist deren FROM zusammen mit der importierten Relation *organization*. Von diesem kartesischen Produkt zieht man das obige Ergebnis ab (NOT EXISTS):



Das Ergebnis dieses Ausdruckes wird sofort auf die Schlüsselattribute der korrelierenden Relation *organization* projiziert (gibt diejenigen Organisationen *O*, für die ein Land mit einer Nichtmitgliedschaft in *O* existiert).

- Das wird jetzt noch von der Menge aller Organisationen abgezogen:



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 is_member
      WHERE is_member.organization=organization.abbreviation
            AND is_member.country=CC));

```

Aufgabe 6 (Optimierung, Aufwandsabschätzung (SQL)) Gegeben sei folgendes Datenbankschema (Auszug aus Mondial)

```

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

```

Formulieren Sie die folgenden Anfragen in SQL:

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

Untersuchen Sie bei der 2. Teilaufgabe den Aufwand (Anzahl Hintergrundspeicherzugriffe, Vergleiche zur Ausführung des Join, Vergleiche zum Test der Fläche). An welchen Stellen wäre ein (Baum- oder anderer) Index nützlich? Untersuchen Sie auch für diesen Fall den Aufwand.

Können Sie Ihre Anfrage optimieren, um damit den Aufwand weiter zu reduzieren?

- a) Welche Länder sind Mitglied in mehr als 60 Organisationen?

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

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

Naheliegendste Lösung: WHERE-Klausel, um vor der Gruppierung zu selektieren.

```
SELECT country
FROM is_member
WHERE (select area
      from country
      where code=is_member.country) > 500000
GROUP BY country
HAVING count(*) > 60
```

Aufwandsbetrachtung: Für jede der 7000 Mitgliedschaften wird überprüft, ob das betreffende Land grösser als 500000 km² ist.

Fall 1: kein Index auf country. Also wird linear die gesamte Tabelle country durchsucht. Aufwand: 7000 · 200 = 1400000 Vergleiche von `code=is_member.country` und 7000 Flächentests.

Fall 2: Baum-Index auf `country.code` [Tafelbild]. Dann wird für jede der Mitgliedschaften über den Baum zielsicher auf das entsprechende `country`-Tupel zugegriffen und die Fläche getestet (7000 Tests) (Dasselbe gilt auch für einen Hash-Index).

Optimierung: der Flächentest wird z.B. für Deutschland 68 mal durchgeführt - für jede Mitgliedschaft einmal. Es wäre effizienter, doch erst nach "country" zu gruppieren, und den Test dann nur für jede Gruppe einmal durchzuführen:

```
SELECT country
FROM is_member
GROUP BY country
HAVING count(*) > 60
and (select area
     from country
     where code=is_member.country) > 500000
```

Damit: 200 Gruppen, bei nested-loop-Join also noch 200·200 = 40000 Vergleiche von `code=is_member.country` und 200 Flächentests, bei Baumindex nur noch 200 Flächentests.

Weitere interne algorithmische Optimierung:

Bei der Auswertung von GROUP BY wird der Inhalt der `is_member`-Tabelle sowieso sortiert. Wenn man nun noch einen Index über `country.code` hat, kann der `area`-Test auch als Merge ablaufen, wobei parallel die sortierten Gruppen und die Blätter des Indexbaumes (was allerdings bei einem Hash-Index nicht gehen würde!) durchlaufen werden

Hinweise:

- Üblicherweise wird automatisch ein Index über alle Schlüssel sowie Fremdschlüssel erzeugt, da diese häufig für Joins und Selektionen benötigt werden.
 - Bei der Auswertung kann ein DBS auch temporär einen Index erstellen.
-
-

Ausserdem:

- Formulieren Sie die Anfragen vom vorigen Blatt auch in SQL, und formulieren Sie die Anfragen von Aufgabe 3 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/>)
Dort finden Sie auch detaillierte Folien sowie ein Skript zu SQL ...