



## Übung zur Vorlesung *Grundlagen: Datenbanken* im WS13/14

Henrik Mühe (muehe@in.tum.de)

<http://www-db.in.tum.de/teaching/ws1314/dbsys/exercises/>

### Blatt Nr. 14

#### Hausaufgabe 1

Für einen Join-Baum  $T$  sei folgende Kostenfunktion gegeben

$$C_{out}(T) = \begin{cases} 0 & \text{falls } T \text{ eine Basisrelation } R_i \text{ ist} \\ |T| + C_{out}(T_1) + C_{out}(T_2) & \text{falls } T = T_1 \bowtie T_2 \end{cases}$$

Die Kardinalität sei dabei

$$|T| = \begin{cases} |R_i| & \text{falls } T \text{ eine Basisrelation } R_i \text{ ist} \\ (\prod_{R_i \in T_1, R_j \in T_2} f_{i,j}) |T_1| |T_2| & \text{falls } T = T_1 \bowtie T_2 \end{cases}$$

Sei  $p_{i,j}$  das Join Prädikat zwischen  $R_i$  und  $R_j$ , dann sei

$$f_{i,j} = \frac{|R_i \bowtie_{p_{i,j}} R_j|}{|R_i \times R_j|}$$

und die Kardinalität eines Join-Resultats ist  $|R_i \bowtie_{p_{i,j}} R_j| = f_{i,j} |R_i| |R_j|$ .

Gegeben sei eine Anfrage über die Relationen  $R_1, R_2, R_3$  und  $R_4$  mit  $|R_1| = 10, |R_2| = 20, |R_3| = 20, |R_4| = 10$ . Die Selektivitäten der Joins seien  $f_{1,2} = 0.01, f_{2,3} = 0.5, f_{3,4} = 0.01$ , alle nicht gegebenen Selektivitäten sind offensichtlich 1 (Warum?). Berechnen Sie den optimalen (niedrigste Kosten) Join-Tree. Als Vereinfachung reicht es, wenn Sie nur Joins mit Prädikat und keine Kreuzprodukte betrachten.

Es ist kein Algorithmus angegeben. Aufgrund der geringen Anzahl von Relationen ist es möglich, die Kosten aller möglichen Join-Bäume zu berechnen und den kostengünstigsten auszuwählen (Bruteforce).

Zunächst gilt es zu überlegen, für welche Join-Bäume die Kosten tatsächlich zu berechnen sind.

Left-Deep:

$$((R_1 \bowtie R_2) \bowtie R_3) \bowtie R_4 \quad (1)$$

$$((R_4 \bowtie R_3) \bowtie R_2) \bowtie R_1 \quad (2)$$

$$((R_3 \bowtie R_2) \bowtie R_1) \bowtie R_4 \quad (3)$$

$$((R_3 \bowtie R_2) \bowtie R_4) \bowtie R_1 \quad (4)$$

Bushy:

$$(R_1 \bowtie R_2) \bowtie (R_3 \bowtie R_4) \quad (5)$$

Alle anderen Left Deep oder Bushy Trees enthalten Kreuzprodukte oder sind im Bezug auf die Kosten äquivalent. Ersteres entsteht, wenn Relationen in einer Reihenfolge gejoin

werden, in der bei einem der Joins kein Prädikt möglich ist, beispielsweise ist dies für den Left-Deep Tree

$$((R_1 \bowtie R_2) \times R_4) \bowtie R_3$$

der Fall. Im Bezug auf die Kosten bei der gegebenen Kostenfunktion äquivalent sind Join-Trees, bei denen die Kinder eines Join Operators vertauscht wurden, etwa

$$(R_1 \bowtie R_2) \bowtie (R_3 \bowtie R_4)$$

und

$$(R_3 \bowtie R_4) \bowtie (R_1 \bowtie R_2).$$

Im Beispiel müssen lediglich die Kosten für die Join-Reihenfolgen 1, 3 und 5 berechnet werden. Dies liegt am Aufbau der Kostenfunktion sowie den symmetrischen Größen der Relationen sowie ihrer Join Selektivitäten.

Die Berechnung von  $C_{out}((R_1 \bowtie R_2) \bowtie (R_3 \bowtie R_4))$  sei hier exemplarisch in epischer Breite ausgeführt (Machen Sie es selber; Sie erkennen äußerst schnell ein Muster und müssen keine derartigen Formel-Konvolute schreiben):

$$\begin{aligned} & C_{out}((R_1 \bowtie R_2) \bowtie (R_3 \bowtie R_4)) \\ = & |(R_1 \bowtie R_2) \bowtie (R_3 \bowtie R_4)| + C_{out}(R_1 \bowtie R_2) + C_{out}(R_3 \bowtie R_4) \\ = & |(R_1 \bowtie R_2) \bowtie (R_3 \bowtie R_4)| + |R_1 \bowtie R_2| + C_{out}(R_1) + C_{out}(R_2) + |R_3 \bowtie R_4| + C_{out}(R_3) + C_{out}(R_4) \\ = & |(R_1 \bowtie R_2) \bowtie (R_3 \bowtie R_4)| + |R_1 \bowtie R_2| + |R_3 \bowtie R_4| \\ = & f_{1,3} * f_{1,4} * f_{2,3} * f_{2,4} * |R_1 \bowtie R_2| * |R_3 \bowtie R_4| + |R_1 \bowtie R_2| + |R_3 \bowtie R_4| \\ = & 0.5 * (0.01 * 10 * 20) * (0.01 * 20 * 10) + (0.01 * 10 * 20) + (0.01 * 20 * 10) \\ = & 2 + 2 + 2 \\ = & 6 \end{aligned}$$

Die Ergebnisse der anderen Relevanten Join-Reihenfolgen sind:

$$\begin{array}{l|l} ((R_1 \bowtie R_2) \bowtie R_3) \bowtie R_4 & 24 \\ ((R_3 \bowtie R_2) \bowtie R_1) \bowtie R_4 & 222 \\ (R_1 \bowtie R_2) \bowtie (R_3 \bowtie R_4) & 6 \end{array}$$

Der Bushy-Tree 5 ist also optimale Join-Tree.

## Hausaufgabe 2

Gegeben sei eine erweiterbare Hashtabelle mit globaler Tiefe  $t$ . Wie viele Verweise zeigen vom Verzeichnis auf einen Behälter mit lokaler Tiefe  $t'$ ?

In dem Verzeichnis einer Hashtabelle mit globaler Tiefe  $t$  werden  $t$  Bits eines Hashwerts für die Identifizierung eines Verzeichniseintrags verwendet. Für einen Behälter mit lokaler Tiefe  $t'$  sind hingegen nur die ersten  $t'$  Bits dieses Bitmusters relevant.

Mit anderen Worten bedeutet dies, dass alle Einträge, die einen Behälter mit lokaler Tiefe  $t'$  referenzieren, in den ersten  $t'$  Bits übereinstimmen. Da alle Bitmuster bis zur Länge  $t$  in dem Directory aufgeführt sind, unterscheiden sich diese Einträge in den letzten  $t - t'$  Bits.  
⇒ Es gibt somit  $2^{t-t'}$  Einträge im Verzeichnis, die auf denselben Behälter mit lokaler Tiefe  $t'$  verweisen.

### Hausaufgabe 3

Gegeben die Relationen:

- *Spieler*(SpielerID, Name, Alter, Team)
- *Herkunft*(Team, Kontinent)
- *Einsatz*(SpielerID, Datum, Ort, Tore)

Wer wurde Weltmeister? Gegen Sie ein SQL Statement an, welches den Namen des Teams bestimmt.

Hinweis: Das Finalspiel ist das einzige Spiel am letzten Tag der WM. Aggregatfunktionen wie z.B. MIN, MAX und COUNT sind auch für Datumswerte definiert.

```
WITH Finale AS (  
  SELECT s.Team, SUM(e.Tore)  
  FROM Spieler s, Einsatz e  
  WHERE e.Datum = (SELECT MAX(e2.Datum) FROM Einsatz e2)  
  AND s.SpielerID = e.SpielerID  
  GROUP BY s.Team  
(  
  SELECT f.Team  
  FROM Finale f  
  WHERE f.Tore = (SELECT MAX(f2.Tore) FROM Finale f2)
```