

## Kapitel 2: Die AGM-Schranke und andere obere Schranken für die Größe von Anfrageergebnissen

Der Begriff "AGM-Schranke" ist benannt nach den Autoren Atserias, Grohe und Marx des Artikels "Size bounds and query plans for relational joins" in Proc. FOCS 2008, pp. 739-748.

### 2.1 Warm-Up: Die $\Delta$ -Anfrage

Betrachte die Anfrage  $Q_\Delta$  aus Beispiel 1.2

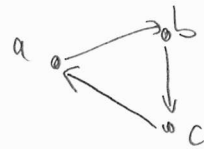
$$Q_\Delta(A, B, C) \leftarrow E(A, B), E(B, C), E(C, A).$$

Sei  $\sigma := \sigma_E = \{E\}$ . Jede  $\sigma$ -DB  $D$  entspricht einem gerichteten Graphen  $G^D = (V^D, E^D)$  mit  $V^D := \text{atom}(D)$ .

Klar: Für jede  $\sigma$ -DB  $D$  ist

$$Q_\Delta(D) = \{ (a, b, c) : (a, b) \in E^D \text{ und } (b, c) \in E^D \text{ und } (c, a) \in E^D \}$$

Skizze:



D.h.:  $Q_\Delta(D)$  besteht aus allen gerichteten Dreiecken in  $G^D$ .

Ziel: Gib eine obere Schranke für die Anzahl  $|Q_\Delta(D)|$  von Tupeln in  $Q_\Delta(D)$  an und finde einen möglichst effizienten Algorithmus, der  $Q_\Delta(D)$  berechnet.

Notation:  $n^D := |\text{dom}(D)|$   
 $N^D := |E^D|$

Offensichtliche obere Schranken für  $|Q_\Delta(D)|$ :

1)  $|Q_\Delta(D)| \leq (n^D)^3$

2)  $|Q_\Delta(D)| \leq (N^D)^2$

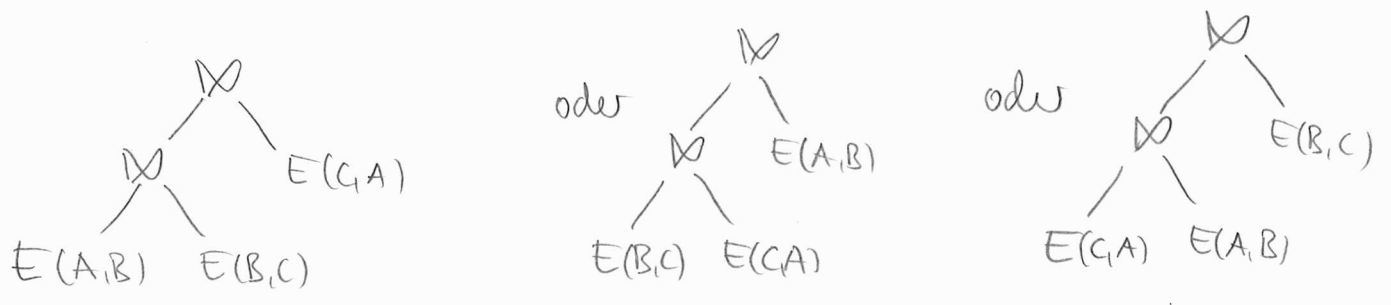
Frage: Geht das besser?

Antwort, die der nächste Satz (Satz  $\Delta$ ) gibt:

Ja:  $|Q_\Delta(D)| \leq 2 \cdot (N^D)^{1,5}$  — und  
 es gibt auch einen Algorithmus, der  $Q_\Delta(D)$   
 mit Laufzeit  $\approx (N^D)^{1,5}$  löst.

Bevor wir uns den Satz (linker Beweis) anschauen,  
 diskutieren wir aber erstmal, wie ein  
 herkömmliches Datenbanksystem die Anfrage  $Q_\Delta$   
 auswertet (siehe Beispiel 1.1 für eine Formulierung  
 der Anfrage in SQL).

Um die Anfrage  $Q_\Delta$  auszuwerten wird ein Datenbanksystem i.d.R. QEPs (query evaluation plans) der folgenden Art betrachtet



Bei jedem dieser QEPs wird als Zwischenergebnis  $Q_2(D)$  berechnet, für

$$Q_2(x,y,z) \leftarrow E(x,y), E(y,z)$$

und dann wird

$$Q_\Delta(x,y,z) \leftarrow Q_2(x,y,z), E(z,x)$$

auf  $D$  ausgewertet.

## Beispiel 2.1

Für jede Zahl  $m \geq 1$  geben wir eine  $\sigma$ -DB  $D_m$  mit

$$N^{D_m} = \Theta(m), \quad |Q_{\Delta}(D_m)| = \Theta(m), \quad \text{aber}$$

$$|Q_2(D_m)| = \Omega(m^2) \quad \text{für} \quad Q_2(X, Y, Z) \leftarrow E(X, Y), E(Y, Z).$$

Die "üblichen" QEPs (query evaluation plans), die von Datenbanksystemen zur Auswertung der Anfrage  $Q_{\Delta}$  erzeugt werden, benötigen zum Auswerten von  $Q_{\Delta}$  auf  $D_m$  daher Zeit  $\Omega((N^{D_m})^2)$ , während der Algorithmus aus dem folgenden Satz  $\Delta$  nur Zeit  $O((N^{D_m})^{1.5})$  benötigt.

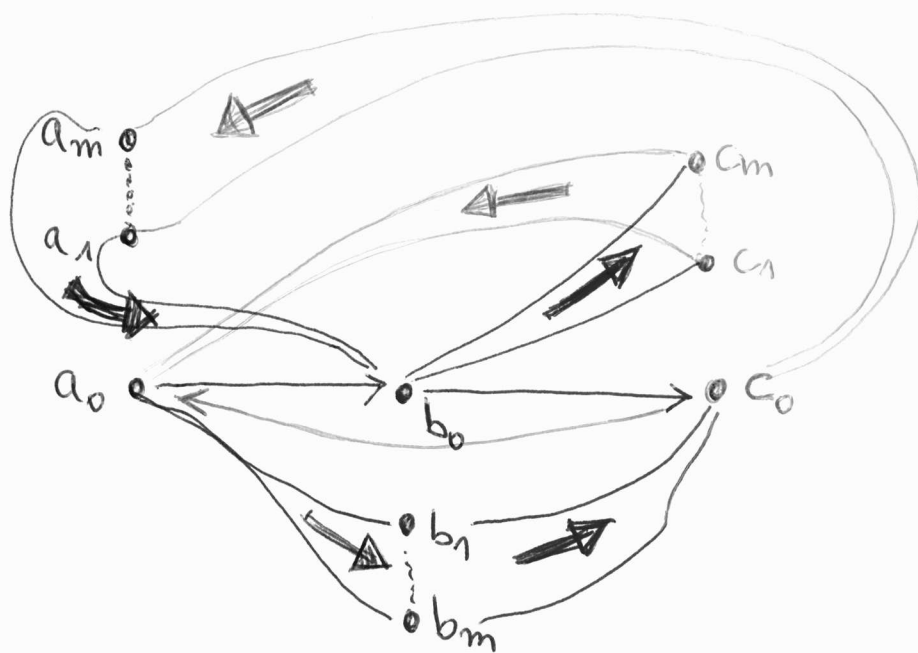
Konstruktion von  $D_m$ :

Seien  $a_0, a_1, \dots, a_m, b_0, b_1, \dots, b_m, c_0, c_1, \dots, c_m$   $3(m+1)$  verschiedene Elemente aus  $\text{dom}$ . Sei

$$E^{D_m} := \left\{ \begin{array}{l} \{ (a_0, b_i) : i \in [0, m] \} \cup \\ \{ (a_i, b_0) : i \in [1, m] \} \cup \\ \{ (b_0, c_i) : i \in [0, m] \} \cup \\ \{ (b_i, c_0) : i \in [1, m] \} \cup \\ \{ (c_0, a_i) : i \in [0, m] \} \cup \\ \{ (c_i, a_0) : i \in [1, m] \} \end{array} \right.$$

klar:  $N^{D_m} := |E^{D_m}| = 3 \cdot (m+1+m) = 6m+3$

Skizze:  $D_m$



Man kann sich leicht davon überzeugen, dass gilt:

$$Q_{\Delta}(D_m) = X \cup \pi_{2,3,1}(X) \cup \pi_{3,1,2}(X) \quad \text{für}$$

$$X = \{ (a_0, b_0, c_i) : i \in [0, m] \} \cup \\ \{ (a_0, b_i, c_0) : i \in [1, m] \} \cup \\ \{ (a_i, b_0, c_0) : i \in [1, m] \}.$$

Es gilt:  $|X| = m+1 + 2m = 3m+1$  und  $|Q_{\Delta}(D_m)| = 3|X| = 9m+3$

Außerdem ist  $(a_i, b_0, c_j) \in Q_2(D_m)$ , f.a.  $i, j \in [0, m]$ ,

also  $|Q_2(D_m)| > m^2$

Satz Δ: Für jede  $\sigma$ -DB  $D$  gilt:

$$|Q_{\Delta}(D)| \leq 2 \cdot (N^D)^{3/2}$$

und es gibt einen Algorithmus, der bei Eingabe von  $D$  die Menge  $Q_{\Delta}(D)$  in Zeit

$O((N^D)^{3/2} \cdot t)$  berechnet, wobei  $t$  die Zeit ist, die wir benötigen, um bei Eingabe von  $(v,w)$  zu testen, ob  $(v,w) \in E^D$  ist.

Beweis: Sei  $D$  eine beliebige  $\sigma$ -DB.

Für jedes  $v \in \text{atom}(D)$  sei

$$\begin{aligned} \text{aus-Grad}^D(v) &:= |\{w : (v,w) \in E^D\}| \\ &= |E^D(v,*)| \end{aligned}$$

$$\text{für } E^D(v,*) := \{w : (v,w) \in E^D\}.$$

Ein Knoten  $v \in \text{atom}(D)$  heißt

- schwer, wenn  $\text{aus-Grad}^D(v) \geq \sqrt{N^D}$  ist
- leicht, wenn  $v$  nicht schwer ist.

Ein Tupel  $t = (a,b,c) \in Q_{\Delta}(D)$  heißt

schwer, wenn  $a$  schwer ist, und es heißt  
 leicht, wenn  $a$  leicht ist.

Für  $x \in \{\text{leicht}, \text{schwer}\}$  sei

$$Q_{\Delta, x}(D) := \{t \in Q_{\Delta}(D) : t \text{ ist } x\}$$

klar:  $Q_{\Delta}(D) = Q_{\Delta, \text{leicht}}(D) \cup Q_{\Delta, \text{schwer}}(D)$

und  $|Q_{\Delta}(D)| = |Q_{\Delta, \text{leicht}}(D)| + |Q_{\Delta, \text{schwer}}(D)|$ .

Behauptung 1:  $|Q_{\Delta, \text{leicht}}(D)| \leq (N^D)^{3/2}$

Beweis: Für jeden leichten Knoten  $a \in \text{adom}(D)$  ist  $\text{aus-Grad}^D(a) < \sqrt{N^D}$ , und es gilt:

$$\begin{aligned} & |\{(b,c) : (a,b,c) \in Q_{\Delta}(D)\}| \\ & \leq |\{(b,c) : (a,b) \in E^D \text{ und } (c,a) \in E^D\}| \\ & \leq \text{aus-Grad}^D(a) \cdot \text{ein-Grad}^D(a), \text{ wobei} \end{aligned}$$

$$\begin{aligned} \text{ein-Grad}^D(a) & := |\{w : (w,a) \in E^D\}| = |E^D(*,a)| \\ \text{für } E^D(*,a) & := \{w : (w,a) \in E^D\}. \end{aligned}$$

Insgesamt gilt:

$$\begin{aligned} |Q_{\Delta, \text{leicht}}(D)| & \leq \sum_{\substack{a \in \text{adom}(D), \\ a \text{ leicht}}} |\{(b,c) : (a,b,c) \in Q_{\Delta}(D)\}| \\ & \leq \sum_{\substack{a \in \text{adom}(D), \\ a \text{ leicht}}} \text{aus-Grad}^D(a) \cdot \text{ein-Grad}^D(a) \\ & \leq \sqrt{N^D} \cdot \sum_{a \in \text{adom}(D)} \text{ein-Grad}^D(a) \\ & \leq \sqrt{N^D} \cdot |E^D| \\ & = \sqrt{N^D} \cdot N^D = (N^D)^{3/2} \end{aligned}$$

□ Beh 1.

Sei  $S^D$  die Menge aller schweren Knoten von  $D$ ,  
d.h.  $S^D := \{v \in \text{dom}(D) : v \text{ ist schwer}\}$ .

20

Behauptung 2:  $|Q_{\Delta, \text{schwer}}(D)| \leq |S^D| \cdot N^D$

Beweis:  $Q_{\Delta, \text{schwer}}(D)$   
 $= \{(a, b, c) : a \in S^D \text{ und } (a, b, c) \in Q_{\Delta}(D)\}$   
 $\subseteq \{(a, b, c) : a \in S^D \text{ und } (b, c) \in E^D\}$   
 $= S^D \times E^D$

Somit:  $|Q_{\Delta, \text{schwer}}(D)| \leq |S^D \times E^D| = |S^D| \cdot |E^D| = |S^D| \cdot N^D$

□ Beh 2

Behauptung 3:  $|S^D| \leq \sqrt{N^D}$

Beweis: Für jedes  $v \in S^D$  ist  $|E^D(v, *)| = \text{aus-Grad}^D(v) \geq \sqrt{N^D}$ .

Es gilt:

$$\begin{aligned} N^D = |E^D| &\geq \left| \bigcup_{v \in S^D} \{(v, w) : w \in E^D(v, *)\} \right| \\ &= \sum_{v \in S^D} |\{(v, w) : w \in E^D(v, *)\}| \\ &= \sum_{v \in S^D} |E^D(v, *)| \\ &\geq \sum_{v \in S^D} \sqrt{N^D} \\ &= |S^D| \cdot \sqrt{N^D} \end{aligned}$$

Somit ist  $|S^D| \leq \frac{N^D}{\sqrt{N^D}} = \sqrt{N^D}$ .

□ Beh 3



Insgesamt gilt:

$$\begin{aligned}
|Q_{\Delta}(D)| &= |Q_{\Delta, \text{leicht}}(D)| + |Q_{\Delta, \text{schwer}}(D)| \\
&\stackrel{\text{Beh 1\&2}}{\leq} (N^D)^{3/2} + |S^D| \cdot N^D \\
&\stackrel{\text{Beh 3}}{\leq} (N^D)^{3/2} + \sqrt{N^D} \cdot N^D \\
&= (N^D)^{3/2} + (N^D)^{3/2} \\
&= 2 \cdot (N^D)^{3/2}
\end{aligned}$$

Unser Algorithmus zur Berechnung von  $Q_{\Delta}(D)$  geht bei Eingabe von  $D$  wie folgt vor:

0) Initialisiere  $Q_{\Delta}(D) := \emptyset$

1) Berechne für jedes  $v \in \text{atom}(D)$  folgendes:

die Mengen  $E^D(v, *) = \{w : (v, w) \in E^D\}$  und  $E^D(*, v) = \{w : (w, v) \in E^D\}$  und

die Zahl  $\text{ans-Grad}^D(v)$ .

Durch geeignetes Sortieren von  $E^D$  geht das in Zeit  $O(N^D \cdot \log N^D) \leq O((N^D)^{3/2})$

2) Berechne die Menge  $S^D$  aller schweren  $v \in \text{atom}(D)$  und die Menge  $L^D := \text{atom}(D) \setminus S^D$  aller leichten  $v \in \text{atom}(D)$ . Unter Verwendung der in 1) gesammelten Infos geht das in Zeit  $O(N^D)$ .

3) Betrachte jedes leichte  $a \in L^D$

Betrachte jedes  $b \in E^D(a, *)$

Betrachte jedes  $c \in E^D(*, a)$

und teste, ob  $(b, c) \in E^D$ .

Wenn ja, füge  $(a, (b, c))$  in  $Q_{\Delta}(D)$  ein

Auf die gleiche Art wie im Beweis von Beh 1 erhalten wir, dass das in Zeit  $O((N^D)^{3/2} \cdot t)$  geht.



4) Betrachte jedes schwere  $a \in S^D$   
 Betrachte jedes  $(b, c) \in E^D$   
 und teste, ob  $(a, b) \in E^D$   
 und  $(c, a) \in E^D$ .



Wenn ja, füge  $(a, b, c)$  in  $Q_\Delta(D)$  ein.

Das geht in Zeit  $|S^D| \cdot N^D \cdot 2t = O((N^D)^{3/2} \cdot t)$ .

Insgesamt berechnen wir so  $Q_\Delta(D)$  in Zeit  $O((N^D)^{3/2} \cdot t)$ .

□ Satz  $\Delta$

Bemerkung  $\Delta$ :

Der Algorithmus aus Satz  $\Delta$  ist im folgenden Sinn Worst-case-optimal:

Es gibt eine Folge von  $\sigma$ -Datenbanken  $D_1, D_2, D_3, \dots$   
 so dass für  $N_i := |E^{D_i}|$  für  $i \in \mathbb{N}_{\geq 1}$  gilt:

$$N_1 < N_2 < N_3 < \dots \text{ und}$$

$$|Q_\Delta(D_i)| = (N_i)^{3/2}.$$

Dazu wähle für  $i \geq 1$  die DB  $D_i$  mit  $E^{D_i} = [1, i] \times [1, i]$ .

Dann ist  $N_i = i^2$  und  $Q_\Delta(D_i) = [1, i]^3$ , also

$$|Q_\Delta(D_i)| = i^3 = (N_i)^{3/2}.$$

## Literaturhinweis:

Der hier präsentierte Beweis von Satz  $\Delta$  ist aus der folgenden Arbeit entnommen:

"Worst-case optimal join algorithms" von  
 Ngo, Pócs, Ré, Rudra,  
 In Proc. PODS 2012, pp. 37-48

Die Beweisidee wird dort Loomis und Whitney zugeschrieben ("An inequality related to the isoperimetric inequality" von Loomis, Whitney. In Bull. Amer. Math. Soc., 55: 361-362, 1949).

# 2.2 Die AGM-Schranke

Die AGM-Schranke wurde zunächst für sog. Join-Anfragen formuliert und später dann auf beliebige konjunktive Anfragen verallgemeinert.

## Definition 2.2

Eine Join-Anfrage (vom Schema  $\sigma$ ) ist eine CQ  $Q$  (vom Schema  $\sigma$ ) der Form

$$Q(\bar{X}_0) \leftarrow R_1(\bar{X}_1), \dots, R_m(\bar{X}_m) \quad (*)$$

für die gilt:

- 1)  $cons(Q) = \emptyset$ ,
- 2) im Kopf der Anfrage kommen alle Variablen vor, die im Rumpf der Anfrage vorkommen
- 3)  $R_1, \dots, R_m$  sind paarweise verschieden, und
- 4) für jedes  $i \in [0, m]$  gilt:  
 $\bar{X}_i$  ist ein Tupel von paarweise verschiedenen Variablen

### Beispiel 2.3

- $Q_{\Delta}(A, B, C) \leftarrow E(A, B), E(B, C), E(C, A)$   
ist keine Join-Anfrage (da Bedingung 3) verletzt wird)
- $Q'_{\Delta}(A, B, C) \leftarrow E_1(A, B), E_2(B, C), E_3(C, A)$   
ist eine Join-Anfrage über dem Schema  
 $\sigma' := \{E_1, E_2, E_3\}$ , das aus drei verschiedenen  
Relationssymbolen der Stelligkeit 2 besteht.
- $Q''_{\Delta}() \leftarrow E_1(A, B), E_2(B, C), E_3(C, A)$   
ist keine Join-Anfrage (da Bedingung 2) verletzt wird)
- Die Anfrage  

$$Q_{(a)}(A) \leftarrow E(A, \text{Sascha Lobo}), E(A, \text{resut öxl})$$
 aus Beispiel 1.2(a) ist keine Join-Anfrage,  
(da sie die Bedingungen 1) und 3) verletzt).

### Bemerkung 2.4

Die in Definition 2.2 eingeführten Join-Anfragen entsprechen in der "benannten Perspektive" der Situation, in der jedes  $R_i$  die Attribute  $\bar{X}_i$  hat, und in der der "natürliche Join"  
 $R_1[\bar{X}_1] \bowtie \dots \bowtie R_m[\bar{X}_m]$  berechnet werden soll.

Die AGM-Schranke liefert eine Antwort auf die folgende Frage:

Gegeben sei eine Join-Anfrage  $Q$  der Form  $\otimes$  und eine DB  $D$  mit  $N_i := N_i^D := |R_i^D|$  f.a.  $i \in [1, m]$ .

Frage: Wie viele Tupel kann es im Anfrageergebnis  $Q(D)$  höchstens geben?

Betrachten wir dazu zunächst einige Spezialfälle:

1) Wenn es ein  $i \in [1, m]$  gibt, s.d.  $\overline{X_0} = \overline{X_i}$  ist, dann ist

$$Q(D) \subseteq R_i^D, \text{ und somit } |Q(D)| \leq N_i$$

2) Wenn  $\overline{X_0} = \overline{X_{i_1}} \cdot \overline{X_{i_2}}$  für  $i_1, i_2 \in [1, m]$  ist, dann ist

$$Q(D) \subseteq R_{i_1}^D \times R_{i_2}^D \text{ und somit } |Q(D)| \leq N_{i_1} \cdot N_{i_2}$$

3) Wenn  $\overline{X_0} = \overline{X_{i_1}} \cdot \dots \cdot \overline{X_{i_k}}$  für  $i_1, \dots, i_k \in [1, m]$  ist,

$$\text{dann ist } Q(D) \subseteq R_{i_1}^D \times \dots \times R_{i_k}^D \text{ und somit } |Q(D)| \leq \prod_{j=1}^k N_{i_j}$$

Letztendlich kommt es dabei nicht auf die Reihenfolge an, in der die einzelnen Variablen in einem Tupel  $\bar{x}_i$  auftreten.

So gilt an Stelle von 1), 2), 3) sogar

1') Wenn es ein  $i \in [1, m]$  gibt, s.d.

$$\{\bar{x}_0\} = \{\bar{x}_i\}, \text{ dann ist } |Q(D)| \leq N_i$$

Hierbei verwenden wir folgende Notation:

Für ein Tupel  $\bar{t} = (t_1, \dots, t_e)$  ist

$$\{\bar{t}\} := \{t_1, \dots, t_e\}.$$

Für eine Liste von Tupeln  $\bar{t}_1, \dots, \bar{t}_k$  ist

$$\{\bar{t}_1, \dots, \bar{t}_k\} := \{\bar{t}_1\} \cup \dots \cup \{\bar{t}_k\}.$$

2') Wenn  $\{\bar{x}_0\} = \{\bar{x}_{i_1}, \bar{x}_{i_2}\}$  für  $i_1, i_2 \in [1, m]$  ist, dann ist  $|Q(D)| \leq N_{i_1} \cdot N_{i_2}$

3') Wenn  $\{\bar{x}_0\} = \{\bar{x}_{i_1}, \dots, \bar{x}_{i_k}\}$  für  $i_1, \dots, i_k \in [1, m]$  ist, dann ist  $|Q(D)| \leq \prod_{j=1}^k N_{i_j}$

Wir formulieren nun 3') von einem anderen Blickwinkel aus:

Definition 2.4 (Anfrage-Graph  $G_Q$ )

Sei  $Q$  eine Join-Anfrage der Form

$$Q(\bar{X}_0) \leftarrow R_1(\bar{X}_1), \dots, R_m(\bar{X}_m)$$

Sei  $n := |\{\bar{X}_1, \dots, \bar{X}_m\}|$  die Anzahl der verschiedenen in  $Q$  vorkommenden Variablen und sei  $A_1, \dots, A_n$  eine Liste all dieser Variablen.

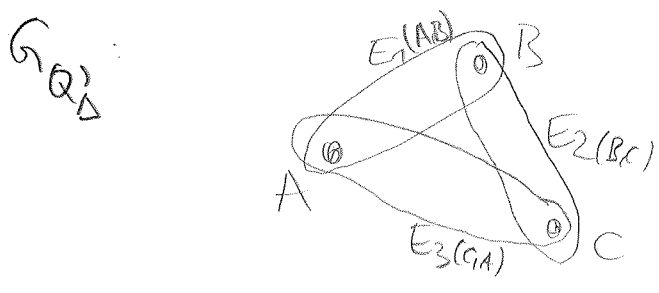
Der zu  $Q$  gehörende Hypergraph  $G_Q$

hat die Knotenmenge  $V_Q := \{A_1, \dots, A_n\} = \text{vars}(Q)$  und für jedes  $i \in [1, m]$  hat  $G_Q$  eine Hyperkante (namens  $R_i(\bar{X}_i)$ ), die aus den Knoten  $\{\bar{X}_i\}$  besteht.

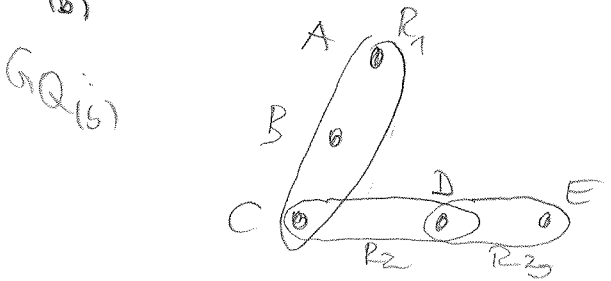
Beachte: Gemäß Bedingung 2) in Def. 2.2 gilt:  $\{\bar{X}_0\} = \{A_1, \dots, A_n\}$

Beispiel 2.5

(a)  $Q_{\Delta}^1(A, B, C) \leftarrow E_1(A, B), E_2(B, C), E_3(C, A)$



(b)  $Q_{(b)}(A, B, C, D, E) \leftarrow R_1(A, B, C), R_2(C, D), R_3(D, E)$





## Definition 2.6 (Kantenüberdeckung)

Sei  $\mathcal{G}$  ein beliebiger Hypergraph mit Knotenmenge  $V$ .

Eine Kantenüberdeckung (engl.: edge cover) von  $\mathcal{G}$  ist eine Menge  $K$  von Hyperkanten von  $\mathcal{G}$ ,

für die gilt:  $\bigcup_{e \in K} e = V$  (d.h.: jeder Knoten

von  $\mathcal{G}$  ist in mindestens einer Hyperkante enthalten, die zu  $K$  gehört).

## Beispiel 2.7

$K := \{E_1(A,B), E_2(B,C)\}$  ist eine Kantenüberdeckung von  $\mathcal{G}_{Q'_\Delta}$  für die Anfrage  $Q'_\Delta$  aus Bsp. 2.5(a)

$K := \{R_1(A,B,C), R_3(D,E)\}$  ist eine Kantenüberdeckung von  $\mathcal{G}_{Q_{(b)}}$  für die Anfrage  $Q_{(b)}$  aus Bsp. 2.5(b)

## Notation:

Wenn wir Join-Anfragen  $Q$  betrachten, sprechen wir von "Kantenüberdeckungen von  $Q$ " und meinen damit Kantenüberdeckungen von  $\mathcal{G}_Q$ .

Die Beobachtung 3') lässt sich also wie folgt formulieren:

3'') Wenn  $i_1, \dots, i_k \in [1, m]$  so gewählt sind, dass  $K := \{ R_{i_1}(\bar{x}_{i_1}), \dots, R_{i_k}(\bar{x}_{i_k}) \}$  eine Kantenüberdeckung von  $Q$  ist, dann gilt:  $|Q(D)| \leq \prod_{j=1}^k N_{i_j}$ .

Eine solche Menge  $K$  können wir auch repräsentieren als eine Abbildung  $x: [1, m] \rightarrow \{0, 1\}$  mit  $x(i) = \begin{cases} 1 & \text{falls } R_i(\bar{x}_i) \in K \\ 0 & \text{falls } R_i(\bar{x}_i) \notin K. \end{cases}$

Dann repräsentiert eine Abbildung  $x: [1, m] \rightarrow \{0, 1\}$  genau dann eine Kantenüberdeckung von  $Q$ , wenn gilt:

Für alle  $j \in [1, m]$  gilt:  $\sum_{\substack{i \in [1, m] \text{ mit} \\ A_j \in \{\bar{x}_i\}}} x(i) \geq 1$  (\*\*)

Beobachtung 3'') lässt sich also auch wie folgt formulieren:

3''') Für jede Abbildung  $x: [1, m] \rightarrow \{0, 1\}$ , die (\*\*) erfüllt, gilt:  $|Q(D)| \leq \prod_{i=1}^m (N_i)^{x(i)}$

Beachte: F. a. reellen Zahlen  $z \neq 0$  gilt:  $z^0 = 1$  und es gilt:  $0^0 = 0$ .

Die AGM-Schranke besagt, dass die Ungleichung  $(***)$  sogar für jede Abbildung  $x: [1, m] \rightarrow \mathbb{Q}_{\geq 0}$ , die  $(**)$  erfüllt, gilt.

Notation:  $\mathbb{Q}_{\geq 0}$  bezeichnet die Menge aller nicht-negativen rationalen Zahlen.

Präzise Formulierung der AGM-Schranke:

Definition 2.8 (Faktionale Kantenüberdeckung)

Sei  $G$  ein beliebiger Hypergraph mit Knotenmenge  $V$  und Hyperkantenmenge  $E$ .

Eine faktionale Kantenüberdeckung (kurz: fKü, engl.: fractional edge cover) von  $G$  ist eine

Abbildung  $x: E \rightarrow \mathbb{Q}_{\geq 0}$ , so dass für jeden

Knoten  $v \in V$  gilt: 
$$\sum_{\substack{e \in E: \\ v \in e}} x(e) \geq 1.$$

Notation:

Wenn wir Join-Anfragen  $Q$  betrachten, sprechen wir von "faktionalen Kantenüberdeckungen von  $Q$ " und meinen damit faktionale Kantenüberdeckungen des Anfrage-Hypergraphen  $G_Q$ . Statt " $x(R_i(\bar{x}_i))$ " schreiben wir auch kurz: " $x(i)$ ".

## Satz 2.9 (Die AGM-Schranke)

Sei  $Q$  eine Join-Anfrage der Form

$$Q(\bar{X}_0) \leftarrow R_1(\bar{X}_1), \dots, R_m(\bar{X}_m).$$

Für jede fraktionale Kantenüberdeckung  $x$  von  $Q$  und jede Datenbank  $D$  gilt:

$$|Q(D)| \leq \prod_{i=1}^m N_i^{x(i)},$$

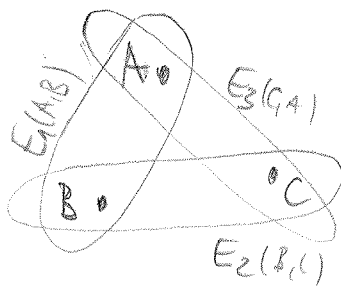
wobei  $N_i := |R_i^D|$  für alle  $i \in [1, m]$  ist.

Bevor wir die AGM-Schranke beweisen, schauen wir uns zunächst einige Beispiele an

### Beispiel 2.10

(a) Betrachte die Anfrage  $Q'_\Delta(A, B, C) \leftarrow E_1(A, B), E_2(B, C), E_3(C, A)$

Skizze von  $G_{Q'_\Delta}$ :



Die Abbildung  $x$  mit  $x(1) = x(2) = x(3) = \frac{1}{2}$  ist eine f.kü von  $Q'_\Delta$ .

Laut AGM-Schranke gilt:  $|Q'_\Delta(D)| \leq \prod_{i=1}^3 N_i^{1/2} = \sqrt{N_1} \cdot \sqrt{N_2} \cdot \sqrt{N_3}$   
f.a. DBen  $D$  und  $N_i := |E_i^D|$  für  $i \in \{1, 2, 3\}$ .

Dies liefert auch eine leicht verbesserte Variante von Satz  $\Delta$ :

Für die Anfrage  $Q_{\Delta}(A,B,C) \leftarrow E(A,B), E(B,C), E(C,A)$  und jede  $\{E\}$ -DB  $D$  gilt:

$$|Q_{\Delta}(D)| \leq (N^D)^{3/2},$$

wobei  $N^D := |E^D|$  ist

Beweis:

Sei  $D$  eine  $\{E\}$ -DB und sei  $D'$  die  $\{E_1, E_2, E_3\}$ -DB mit  $E_i^{D'} := E^D$  f.a.  $i \in [1,3]$ . Dann gilt:

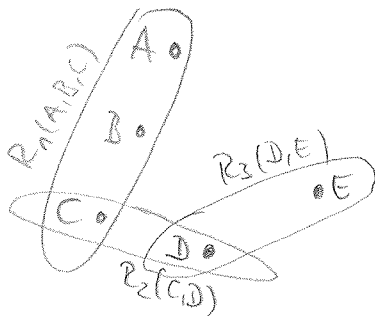
$$Q_{\Delta}(D) = Q'_{\Delta}(D'). \quad \text{Gemäß AGM-Schranke folgt:}$$

$$|Q_{\Delta}(D)| = |Q'_{\Delta}(D')| \leq \sqrt{N^{D'}} \cdot \sqrt{N^{D'}} \cdot \sqrt{N^{D'}} = (N^D)^{3/2}.$$

(b) Betrachte die Anfrage  $Q_{(b)}$  aus Beispiel 2.5:

$$Q_{(b)}(A,B,C,D,E) \leftarrow R_1(A,B,C), R_2(C,D), R_3(D,E).$$

$G_{Q_{(b)}}$ :

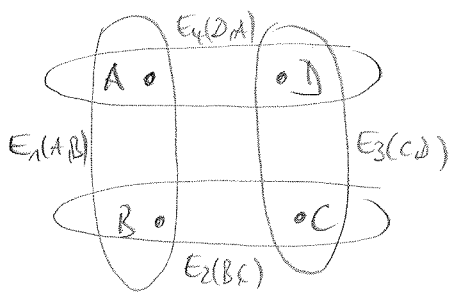


Für jede f.kü  $x$  von  $Q_{(b)}$  gilt:  $x(1) \geq 1$  und  $x(3) \geq 1$ , da der Knoten  $A$  (bzw.  $E$ ) nur in einer einzigen Hyperkante, nämlich  $R_1(A,B,C)$  (bzw.  $R_3(D,E)$ ) vorkommt.

(c) Betrachte die Anfrage

$$Q_4(A, B, C, D) \leftarrow E_1(A, B), E_2(B, C), E_3(C, D), E_4(D, A)$$

$Q_4$ :



Einige fkt. von  $Q_4$ :

	1	2	3	4
$x^{(1)}$	1	0	1	0
$x^{(2)}$	0	1	0	1
$x^{(3)}$	$1/2$	$1/2$	$1/2$	$1/2$
$x^{(4)}$	$1/3$	$2/3$	$1/3$	$2/3$

(d) Sei  $k \in \mathbb{N}$  mit  $k \geq 2$ . Betrachte die  
Loomis-Whitney Join-Anfrage  $LW_k$  mit

$$LW_k(x_1, \dots, x_k) \leftarrow R_1(x_2, \dots, x_k), R_2(x_1, x_3, \dots, x_k), \dots, R_k(x_1, \dots, x_{k-1})$$

dh.: für jedes  $i \in [1, k]$  ein Atom der Form  
 $R_i(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_k)$

Eine fkt. von  $LW_k$  ist die Abbildung  $x: [1, k] \rightarrow \mathbb{Q}_{\geq 0}$  mit  
 $x(i) := 1/(k-1)$  f.a.  $i \in [1, k]$ .

Die AGM-Schranke liefert:  $|LW_k(D)| \leq \prod_{i=1}^k \sqrt[k-1]{N_i}$ ,

wobei  $N_i := |R_i^D|$  f.a.  $i \in [1, k]$  ist  
Insbes.: Falls  $N_1 = \dots = N_k = N$ , so ist  $|LW_k(D)| \leq N^{k/(k-1)} = N \cdot \sqrt[k-1]{N}$ .

Betrachten wir nun folgendes Szenario:

Gegeben sei eine Join-Anfrage

$$Q(\bar{X}_0) \leftarrow R_1(\bar{X}_1), \dots, R_m(\bar{X}_m)$$

mit  $\text{vars}(Q) = \{A_1, \dots, A_n\}$ .

Anßerdem seien Zahlen  $N_1, \dots, N_m$  bekannt, von denen wir wissen, dass die Datenbank  $D$ , auf der  $Q$  ausgewertet wird in der Relation  $R_i$  genau  $N_i$  Tupel enthält (d.h.  $N_i = |R_i^D|$ ), f.a.  $i \in [1, m]$

Die AGM-Schranke besagt, dass für jede fkt  $x$  von  $Q$  gilt:

$$|Q(D)| \leq \prod_{i=1}^m N_i^{x(i)}$$

Frage: Wie können wir eine fkt  $x$  finden, für die  $\prod_{i=1}^m N_i^{x(i)}$  möglichst klein ist?

Dazu lösen wir folgendes Optimierungsproblem:

Statt  $x(i)$  schreibe  $x_i$  für  $i \in [1, m]$  und setze  $x := \begin{pmatrix} x_1 \\ \vdots \\ x_m \end{pmatrix}$ .

Ziel: minimiere  $\prod_{i=1}^m N_i^{x_i}$   
 unter der Bedingung, dass Folgendes gilt:

(1) f.a.  $j \in [1, m]$  ist  $\sum_{\substack{i \in [1, m]: \\ A_j \in \bar{X}_i}} x_i \geq 1$   
 und  
 (2) f.a.  $i \in [1, m]$  ist  $x_i \geq 0$ ,  $x_i \in \mathbb{Q}$

Beachte: •  $\prod_{i=1}^m N_i^{x_i}$  ist minimal

$(\Rightarrow)$   $\log\left(\prod_{i=1}^m N_i^{x_i}\right)$  ist minimal

$$= \sum_{i=1}^m \log(N_i^{x_i}) = \sum_{i=1}^m x_i \cdot \log N_i = \sum_{i=1}^m x_i \cdot c_i$$

mit  $c_i := \log N_i$  f.a.  $i \in [1, m]$ .

Als Zielfunktion für unser Minimierungsproblem können wir statt  $\prod_{i=1}^m N_i^{x_i}$  also auch  $c^T x$  verwenden.

• Die Bedingung (1) lässt sich darstellen als

" $Ax \geq b$ " für die Matrix  $A_{\mathbb{Q}} := A := (a_{ji})_{\substack{j \in [1, m], \\ i \in [1, m]}}$  mit

$$a_{ji} := \begin{cases} 1 & \text{falls } A_{ji} \in \{\bar{x}_i\} \\ 0 & \text{sonst} \end{cases}$$

und den Vektor  $b = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}$

Unser Ziel ist also:

minimiere  $c^T x$  (mit  $c_i = \log N_i$   
f.a.  $i \in [1, m]$ )

unter der Bedingung, dass

$$Ax \geq b, \quad \left( \text{mit } A = (a_{ji}) \text{ und } a_{ji} := \begin{cases} 1 & \text{falls } A_{ji} \in \{\bar{x}_i\} \\ 0 & \text{sonst} \end{cases} \right)$$

$$x \geq 0, \quad x \in \mathbb{R}^m$$

Dies ist ein lineares Optimierungsproblem, für das wir eine rationale Lösung suchen.



Dazu stehen verschiedene Lösungsverfahren zur Verfügung:

- der Simplex-Algorithmus (exponentiell im worst-case, aber in der Praxis i.d.R. effizient)
- das Innere-Punkte-Verfahren (Karmarkar-Verfahren (1984) — läuft in Polynomialzeit)
- die Ellipsoid-Methode (Khachiyan (1979) — läuft in Polynomialzeit)

Beachte:

Ist  $x$  eine Lösung des Optimierungsproblems OPT, so gilt für jede DB  $D$  mit  $N_i = |\mathbb{R}_i^D|$  f.a.  $i \in \{1, m\}$ , dass  $|Q(D)| \leq 2^{c^T x}$  ist.

Denn:

$$2^{c^T x} = 2^{\sum_{i=1}^m c_i x_i} = 2^{\sum_{i=1}^m x_i \cdot \log N_i} = 2^{\sum_{i=1}^m \log(N_i^{x_i})}$$

$$= 2^{\log\left(\prod_{i=1}^m N_i^{x_i}\right)} = \prod_{i=1}^m N_i^{x_i}$$

## Beweis der AGM-Schranke

Zum Beweis benötigen wir einige Grundlagen aus der Wahrscheinlichkeitsrechnung und der Informationstheorie, die wir hier zunächst bereitstellen.

### Wahrscheinlichkeitsräume

Definition: Ein endlicher Wahrscheinlichkeitsraum  $(\Omega, P)$  besteht aus einer endlichen, nicht-leeren Menge  $\Omega$  von Ergebnissen bzw. Elementarereignissen, denen Wahrscheinlichkeiten  $P(\omega) = p_\omega \in \mathbb{R}$  für jedes  $\omega \in \Omega$  zugeordnet sind, so dass gilt:

$$0 \leq p_\omega \leq 1, \text{ für jedes } \omega \in \Omega, \text{ und } \sum_{\omega \in \Omega} p_\omega = 1.$$

### Ereignisse

Definition: Ein Ereignis ist eine Menge von Ergebnissen, d.h. eine Teilmenge von  $\Omega$ . Die Wahrscheinlichkeit eines Ereignisses  $A \subseteq \Omega$  ist definiert als

$$P(A) := \sum_{\omega \in A} P(\omega).$$

Wir schreiben  $\bar{A}$  um das Komplement von  $A$  zu bezeichnen, d.h.  $\bar{A} = \Omega \setminus A$

Regeln zum Rechnen mit Wahrscheinlichkeiten:

Für jeden endlichen Wahrscheinlichkeitsraum  $(\Omega, \mathcal{P})$  und für alle Ereignisse  $A$  und  $B$  gilt:

- $P(\bar{A}) = 1 - P(A)$
- $P(A \cup B) = P(A) + P(B) - P(A \cap B)$

Für  $s \in \mathbb{N}_{\geq 1}$  und beliebige Ereignisse  $A_1, \dots, A_s \in \mathcal{P}$

gilt:  $P(A_1 \cup \dots \cup A_s) \leq \sum_{i=1}^s P(A_i)$

Falls  $A_i \cap A_j = \emptyset$  für alle  $i, j \in [1, s]$  mit  $i \neq j$ , so gilt

sogar:  $P(A_1 \cup \dots \cup A_s) = \sum_{i=1}^s P(A_i)$ .

## Zufallsvariablen

Definition:

Sei  $(\Omega, \mathcal{P})$  ein endlicher Wahrscheinlichkeitsraum.

Eine Zufallsvariable ist eine Funktion  $Y: \Omega \rightarrow M$ ,

für eine beliebige nicht-leere Menge  $M$ .

Für jedes  $a \in M$  definieren wir die Wahrscheinlichkeit dafür, dass die Zufallsvariable  $Y$  den Wert  $a$  annimmt, durch

$P(Y=a) := P(A)$  für das Ereignis  $A := \{\omega \in \Omega : Y(\omega) = a\}$ .

### Beispiel 2.11

Betrachte eine Join-Anfrage  $Q$  der Form

$$Q(\bar{X}_0) \leftarrow R_1(\bar{X}_1), \dots, R_m(\bar{X}_m)$$

mit  $\bar{X}_0 = A_1, \dots, A_m$ .

Sei  $D$  eine Datenbank vom Schema  $\{R_1, \dots, R_m\}$ .

Wir betrachten das "Zufallsexperiment", bei dem wir zufällig, gleichverteilt ein beliebiges Tupel aus dem Anfrageergebnis  $Q(D)$  ziehen.

Dies entspricht dem endlichen Wahrscheinlichkeitsraum

$$(\Omega, \mathcal{P}) \text{ mit } \Omega := Q(D) \text{ und } P(\omega) := \frac{1}{|Q(D)|}$$

f.a.  $\omega \in Q(D)$ .

Sei  $M := \text{atom}(D)$ .

Für jedes  $j \in [1, m]$  sei  $Y_j: \Omega \rightarrow M$  die

Zufallsvariable, die jedem Tupel  $t = (t_1, \dots, t_m) \in Q(D)$

den Wert  $t_j$  zuordnet — d.h.  $Y_j(t) = \pi_j(t)$  f.a.  $t \in Q(D)$ .

Für jedes feste Element  $a \in \text{atom}(D)$  ist dann

$P(Y_j = a)$  die Wahrscheinlichkeit dafür, dass ein

zufällig, gleichverteilt aus  $Q(D)$  gewähltes Tupel in der  $j$ -ten Spalte den Eintrag  $a$  hat. D.h.

$$P(Y_j = a) = \frac{|\{t \in Q(D) : \pi_j(t) = a\}|}{|Q(D)|}.$$

# Informationsgehalt und Entropie

## Definition:

Sei  $(\Omega, P)$  ein endlicher Wahrscheinlichkeitsraum

Der Informationsgehalt eines Elementarereignisses

$w \in \Omega$  ist definiert als

$$\log\left(\frac{1}{P_w}\right) \quad (= -\log(P_w))$$

Die Entropie von  $(\Omega, P)$  ist definiert als

$$H((\Omega, P)) := \sum_{w \in \Omega} P_w \cdot \log\left(\frac{1}{P_w}\right) \quad (= -\sum_{w \in \Omega} P_w \cdot \log(P_w))$$

Konvention für den Fall, dass  $P_w = 0$  ist:

$$0 \cdot \log\left(\frac{1}{0}\right) := -0 \cdot \log(0) := 0$$

## Bemerkung:

Man kann beweisen, dass Folgendes gilt:

$$H((\Omega, P)) \leq \left[ \begin{array}{l} \text{durchschnittliche Länge des Bitstrings,} \\ \text{der ein gemäß } P \text{ zufällig aus } \Omega \\ \text{gewähltes Elementarereignis kodiert} \\ \text{(bei optimaler Kodierung)} \end{array} \right] \leq H((\Omega, P)) + 1$$

Wegen  $H((\Omega, P)) \stackrel{\text{Def}}{=} \sum_{w \in \Omega} P_w \cdot \log\left(\frac{1}{P_w}\right)$  gilt daher:

$$\left[ \begin{array}{l} \text{Informationsgehalt} \\ \text{von } w \in \Omega \end{array} \right] \stackrel{\text{Def}}{=} \log\left(\frac{1}{P_w}\right) \approx \left[ \begin{array}{l} \text{Länge des Bitstrings,} \\ \text{der } w \text{ kodiert} \\ \text{(bei optimaler Kodierung)} \end{array} \right]$$

Definition (Entropie einer Zufallsvariablen)

Sei  $(\Omega, P)$  ein endlicher Wahrscheinlichkeitsraum,  
Sei  $M$  eine beliebige nicht-leere Menge und  
Sei  $Y: \Omega \rightarrow M$  eine Zufallsvariable.

Die Entropie von  $Y$  ist definiert als

$$H(Y) := \sum_{a \in M} P(Y=a) \cdot \log\left(\frac{1}{P(Y=a)}\right).$$

Bemerkung:

Wir können uns  $H(Y)$  vorstellen als ein Maß für die  
"Unsicherheit" bzw. "Unbestimmtheit" von  $Y$ .

Falls es ein  $b \in M$  gibt, so das  $P(Y=b) = 1$ ,  
so wissen wir, dass die Zufallsvariable, unabhängig  
vom Ausgang des Zufallsexperiments, stets den  
Wert  $b$  annimmt — die "Unsicherheit" bzw.  
"Unbestimmtheit" von  $Y$  ist also 0.

Die spiegelt sich in der Entropie wider, da

$$\begin{aligned} H(Y) &\stackrel{\text{Def}}{=} \sum_{a \in M} P(Y=a) \cdot \log\left(\frac{1}{P(Y=a)}\right) = P(Y=b) \cdot \log\left(\frac{1}{P(Y=b)}\right) \\ &= 1 \cdot \log\left(\frac{1}{1}\right) = 0 \end{aligned}$$

ist.

Falls jedoch für jedes  $b \in M$  gilt, dass  $P(Y=b) = \frac{1}{|M|}$ ,  
so ist die "Unsicherheit / Unbestimmtheit" von  $Y$  größtmöglich,  
und es gilt:  $H(Y) = \sum_{a \in M} \frac{1}{|M|} \cdot \log(|M|) = \log(|M|)$

# Wichtige Eigenschaften der Entropie:

Sei  $(\Omega, \mathcal{P})$  ein endlicher Wahrscheinlichkeitsraum

(1) Für jede Zufallsvariable  $Y: \Omega \rightarrow M$  mit endlichem  $M$  gilt:

$$0 \leq H(Y) \leq \log(|M|)$$

Anßerdem gilt:

•  $H(Y) = 0 \iff \exists b \in M$  s.d.  $P(Y=b) = 1$

•  $H(Y) = \log(|M|) \iff \forall a, b \in M$  ist  $P(Y=b) = \frac{1}{|M|}$ .

(2) Seien  $Y_1: \Omega \rightarrow M_1$  und  $Y_2: \Omega \rightarrow M_2$  zwei Zufallsvariablen.

Die zusammengesetzte Zufallsvariable  $(Y_1, Y_2)$  ist die Zufallsvariable  $Y: \Omega \rightarrow M_1 \times M_2$  mit

$$Y(\omega) := (Y_1(\omega), Y_2(\omega)) \quad \forall \omega \in \Omega.$$

• Die gemeinsame Entropie von  $Y_1$  und  $Y_2$  ist die Entropie von  $(Y_1, Y_2)$ . D.h.:

$$H(Y_1, Y_2) = \sum_{\substack{a_1 \in M_1 \\ a_2 \in M_2}} P(Y_1=a_1, Y_2=a_2) \cdot \log\left(\frac{1}{P(Y_1=a_1, Y_2=a_2)}\right)$$

↑  
Komma  $\hat{=}$  "und"

• Für  $a_1 \in M_1$  mit  $P(Y_1=a_1) \neq 0$  ist die bedingte Entropie von  $Y_2$  unter der Voraussetzung, dass  $Y_1=a_1$  ist,

$$H(Y_2 | Y_1=a_1) := \sum_{a_2 \in M_2} P(Y_2=a_2 | Y_1=a_1) \cdot \log\left(\frac{1}{P(Y_2=a_2 | Y_1=a_1)}\right)$$

Hierbei ist  $P(Y_2=a_2 | Y_1=a_1)$  die bedingte Wahrscheinlichkeit dafür, dass  $Y_2$  den Wert  $a_2$  annimmt unter der Voraussetzung, dass  $Y_1$  den Wert  $a_1$  angenommen hat.

$$\text{D.h.: } P(Y_2=a_2 | Y_1=a_1) := \frac{P(Y_2=a_2, Y_1=a_1)}{P(Y_1=a_1)}$$

• Die bedingte Entropie von  $Y_2$  bei gegebenem  $Y_1$  ist

$$\begin{aligned} H(Y_2 | Y_1) &:= \sum_{a_1 \in M_1} P(Y_1=a_1) \cdot H(Y_2 | Y_1=a_1) \\ &= \sum_{a_1 \in M_1} \left( P(Y_1=a_1) \cdot \sum_{a_2 \in M_2} P(Y_2=a_2 | Y_1=a_1) \cdot \log \left( \frac{1}{P(Y_2=a_2 | Y_1=a_1)} \right) \right) \end{aligned}$$

• Es gilt: 
$$H(Y_1, Y_2) = H(Y_1) + H(Y_2 | Y_1) \quad \textcircled{1}$$

Beweis:

$$\begin{aligned} H(Y_1, Y_2) &\stackrel{\text{Def}}{=} \sum_{\substack{a_1 \in M_1 \\ a_2 \in M_2}} \overbrace{P(Y_1=a_1, Y_2=a_2)}^{= P(Y_1=a_1) \cdot P(Y_2=a_2 | Y_1=a_1)} \cdot \log \left( \frac{1}{P(Y_1=a_1, Y_2=a_2)} \right) \\ &= \sum_{a_1 \in M_1} \left( P(Y_1=a_1) \cdot \sum_{a_2 \in M_2} P(Y_2=a_2 | Y_1=a_1) \cdot \left( \log \left( \frac{1}{P(Y_1=a_1)} \right) + \log \left( \frac{1}{P(Y_2=a_2 | Y_1=a_1)} \right) \right) \right) \\ &= \sum_{a_1 \in M_1} P(Y_1=a_1) \cdot \log \left( \frac{1}{P(Y_1=a_1)} \right) \cdot \left( \sum_{a_2 \in M_2} P(Y_2=a_2 | Y_1=a_1) \right) = 1 \\ &\quad + \sum_{a_1 \in M_1} P(Y_1=a_1) \cdot \left( \sum_{a_2 \in M_2} P(Y_2=a_2 | Y_1=a_1) \cdot \log \left( \frac{1}{P(Y_2=a_2 | Y_1=a_1)} \right) \right) \\ &= H(Y_1) + H(Y_2 | Y_1) \quad \square \end{aligned}$$



• Es gilt: 
$$H(Y_2 | Y_1) \leq H(Y_2) \quad (2)$$

Intuition dazu: Die "Unsicherheit von  $Y_2$ " kann durch die zusätzliche Information hinsichtlich  $Y_1$  höchstens abnehmen, aber nicht größer werden.

Beweis von (2): Übung!

Hinweis: Verwende die so genannte Jensensche Ungleichung, die Folgendes besagt:

Für jede konkave Funktion  $f$  und für Zahlen  $\lambda_1, \dots, \lambda_n \in \mathbb{R}_{\geq 0}$  mit  $\sum_{i=1}^n \lambda_i = 1$  gilt

$$f\left(\sum_{i=1}^n \lambda_i x_i\right) \geq \sum_{i=1}^n \lambda_i f(x_i) \quad \text{für alle } x_1, \dots, x_n \in \mathbb{R}$$

Eine Funktion  $f: \mathbb{R} \rightarrow \mathbb{R}$  heißt konkav, wenn für alle  $x, y \in \mathbb{R}$  und alle  $\lambda \in \mathbb{R}$  mit  $0 \leq \lambda \leq 1$  gilt:

$$f(\lambda x + (1-\lambda)y) \geq \lambda f(x) + (1-\lambda)f(y)$$

Insbesondere ist die Funktion  $f(x) := \log x$  konkav

(3) Die Begriffe aus (2) können für zusammengesetzte Zufallsvariablen  $(Y_1, \dots, Y_n)$  verallgemeinert werden:  
Sei  $n \geq 2$ . Für jedes  $j \in \{n\}$  sei  $Y_j: \Omega \rightarrow M_j$  eine Zufallsvariable.

Die zusammengesetzte Zufallsvariable  $(Y_1, \dots, Y_n)$  ist die Zufallsvariable  $Y: \Omega \rightarrow M_1 \times \dots \times M_n$  mit

$$Y(\omega) := (Y_1(\omega), \dots, Y_n(\omega)) \quad \text{f.a. } \omega \in \Omega.$$

• Die gemeinsame Entropie von  $Y_1, \dots, Y_n$  ist

$$H(Y_1, \dots, Y_n) := \sum_{(a_1, \dots, a_n) \in M_1 \times \dots \times M_n} P\left(\bigwedge_{j=1}^n Y_j = a_j\right) \cdot \log\left(\frac{1}{P\left(\bigwedge_{j=1}^n Y_j = a_j\right)}\right)$$

• Analog zu (1) gilt

$$H(Y_1, \dots, Y_n) = H(Y_1) + H(Y_2 | Y_1) + H(Y_3 | Y_1, Y_2) + \dots + H(Y_n | Y_1, \dots, Y_{n-1}) \quad (1)$$

• Analog zu (2) gilt f.a.  $J \subseteq J' \subseteq [n]$  und f.a.  $k \in [n]$ :

$$H(Y_k | (Y_j)_{j \in J'}) \leq H(Y_k | (Y_j)_{j \in J}) \quad (2)$$

Das zentrale Werkzeug, das wir zum Beweis der AAM-Schranke verwenden werden, ist

Shearer's Lemma: Sei  $(\Omega, P)$  ein endlicher Wahrscheinlichkeitsraum

Sei  $n \in \mathbb{N}_{\geq 1}$  und sei  $I := \{1, \dots, n\}$ . Für jedes  $j \in I$  sei  $Y_j: \Omega \rightarrow M_j$  eine Zufallsvariable.

Für jedes  $J \subseteq I$  sei  $Y_J = (Y_j)_{j \in J}$ .

Sei  $\ell \in \mathbb{N}_{\geq 1}$ , und seien  $J_1, \dots, J_\ell$  Teilmengen von  $I$ , so dass für eine Zahl  $q \in \mathbb{N}$  gilt: jedes  $j \in I$  kommt in mindestens  $q$  der Mengen  $J_k$  vor.

Dann gilt:  $H(Y_1, \dots, Y_n) \leq \frac{1}{q} \cdot \sum_{k=1}^{\ell} H(Y_{J_k})$

Beweis:

Gemäß (1) gilt für jedes  $J \subseteq I$ , dass

$$H(Y_J) = \sum_{j \in J} H(Y_j | (Y_i)_{i \in J \text{ mit } i < j}) \quad *$$

Somit gilt:

$$\sum_{k=1}^{\ell} H(Y_{J_k})$$

$$\stackrel{=}{=} \sum_{k=1}^{\ell} \sum_{j \in J_k} H(Y_j | (Y_i)_{i \in J_k \text{ mit } i < j}) \quad *$$

$$\stackrel{\geq}{\geq} \sum_{k=1}^{\ell} \sum_{j \in J_k} H(Y_j | (Y_i)_{i \in I \text{ mit } i < j}) \quad (2)$$

$$\stackrel{\geq}{\geq} q \cdot \sum_{j \in I} H(Y_j | (Y_i)_{i \in I \text{ mit } i < j}) \quad = \quad q \cdot H(\underbrace{Y_I}_m) = (Y_1, \dots, Y_m)$$

den laut Voraussetzung gilt:  
 jedes  $j \in I$  kommt in  
 mindestens  $q$  der Mengen  $J_1, \dots, J_k$  vor.

□

Beweis der AGM-Schranke, dh Beweis von Satz 2.9:

Sei  $Q$  eine Join-Anfrage der Form  $Q(\bar{X}_0) \leftarrow R_1(\bar{X}_1), \dots, R_m(\bar{X}_m)$  und sei  $\bar{X}_0 = A_1, \dots, A_m$ .

Sei  $x$  eine fraktionale Kantenüberdeckung von  $Q$ .

Sei  $D$  eine Datenbank vom Schema  $\{R_1, \dots, R_m\}$ , und für jedes  $i \in [1, m]$  sei  $N_i := |R_i^D|$ .

Zu zeigen:  $|Q(D)| \leq \prod_{i=1}^m N_i^{x(i)}$ .

Falls  $Q(D) = \emptyset$  ist, so gilt dies offensichtlich.

Falls  $Q(D) \neq \emptyset$  ist, so betrachte das Zufallsexperiment aus Beispiel 2.11, bei dem zufällig, gleichverteilt ein beliebiges Tupel aus  $Q(D)$  gewählt wird. D.h.: wir betrachten den endlichen Wahrscheinlichkeitsraum  $(\Omega, P)$  mit  $\Omega := Q(D)$  und  $P(t) := \frac{1}{|Q(D)|}$  für jedes  $t \in Q(D)$ .

Sei  $M := \text{adom}(D)$ , und für jedes  $j \in [1, m]$  sei  $Y_j: \Omega \rightarrow M$  die Zufallsvariable mit  $Y_j(t) := \pi_j(t)$  für jedes  $t \in Q(D)$  (dh:  $Y_j$  ordnet jedem Tupel in  $Q(D)$  seinen Eintrag in der  $j$ -ten Spalte zu).

Beachte: Dann ist  $Y := (Y_1, \dots, Y_m)$  die Zufallsvariable  $Y: Q(D) \rightarrow \text{adom}(D)^m$  mit  $Y(t) = t$  f.a.  $t \in Q(D)$ .

Insbes. gilt für jedes  $t \in \text{adom}(D)^m$ , dass

$$P(Y=t) = \begin{cases} 0 & \text{falls } t \notin Q(D) \\ \frac{1}{|Q(D)|} & \text{falls } t \in Q(D). \end{cases}$$

Daher ist

$$H(Y) = \sum_{t \in \text{adom}(D)^m} P(Y=t) \log\left(\frac{1}{P(Y=t)}\right) = \sum_{t \in Q(D)} P(Y=t) \cdot \log\left(\frac{1}{P(Y=t)}\right) = \log(|Q(D)|) \quad \textcircled{1}$$

Laut Voraussetzung ist  $x$  eine fraktionale Kantenüberdeckung von  $Q$ , d.h.:  $x(1), \dots, x(m) \in \mathbb{Q}_{\geq 0}$  und für jedes  $j \in [1, m]$  gilt

$$\sum_{\substack{i \in [1, m] \text{ mit} \\ A_j \in \bar{X}_i}} x(i) \geq 1 \quad (2)$$

Wegen  $x(1), \dots, x(m) \in \mathbb{Q}_{\geq 0}$  gibt es ein  $q \in \mathbb{N}_{\geq 1}$  und  $p_1, \dots, p_m \in \mathbb{N}$ , s.d.  $x(i) = \frac{p_i}{q}$  f.a.  $i \in [1, m]$  ist.

Sei  $I := \{1, \dots, m\}$ . Für jedes  $i \in [1, m]$  sei

$$J^{(i)} := \{j \in [1, m] : A_j \in \bar{X}_i\}, \text{ und sei}$$

$J_1^{(i)}, \dots, J_{p_i}^{(i)}$  eine Liste, die aus  $p_i$  Kopien von  $J^{(i)}$  besteht

Dann ist  $J_1^{(1)}, \dots, J_{p_1}^{(1)}, \dots, J_1^{(m)}, \dots, J_{p_m}^{(m)}$  eine Liste von

$l = p_1 + \dots + p_m$  Teilmengen von  $I$ , für die gilt:

Jedes  $j \in I$  kommt in mindestens  $q$  dieser Mengen vor,

denn: Für jedes  $j \in I = [1, m]$  gilt:

Anzahl der Mengen in  $(J_1^{(i)}, \dots, J_{p_i}^{(i)})_{i \in [1, m]}$ , in denen  $j$  vorkommt

$$= \sum_{\substack{i \in [1, m] \text{ mit} \\ j \in J^{(i)}}} p_i$$

$$= \sum_{\substack{i \in [1, m] \text{ mit} \\ A_j \in \bar{X}_i}} q \cdot \frac{p_i}{q} = q \cdot \sum_{\substack{i \in [1, m] \text{ mit} \\ A_j \in \bar{X}_i}} x(i) \geq q \quad (2)$$

Shearer's Lemma liefert:  $H(Y) \leq \frac{1}{q} \cdot \sum_{\substack{i \in [1, m], \\ k \in [1, p_i]}} H(Y_{J_k^{(i)}}), \quad (3)$

wo sei  $Y_{j^{(i)}} = (Y_j)_{j \in J_k^{(i)}} = (Y_j)_{A_j \in \{X_i\}}$  ist <sup>50</sup>

Beachte:  $Y_{j^{(i)}} := (Y_j)_{A_j \in \{X_i\}}$  ist eine Zufallsvariable

$Y_{j^{(i)}} : \mathcal{Q}(D) \rightarrow \text{adom}(D)^{\text{ar}(R_i)}$ , für die für jedes

$t \in \mathcal{Q}(D)$  gilt:  $Y_{j^{(i)}}(t) \in R_i^D$

(da  $Q$  von der Form  $Q(A_1, \dots, A_n) \leftarrow R_1(X_1), \dots, R_m(X_m)$  ist).

Daher ist

$$H(Y_{j^{(i)}}) \leq \log(|R_i^D|) = \log N_i. \quad (4)$$

(da  $H(Y) = \log(|M|)$  für jede Zufallsvariable  $Y: \mathcal{R} \rightarrow M$  gilt).

Insgesamt erhalten wir:

$$\log(|\mathcal{Q}(D)|) \stackrel{(1)}{=} H(Y) \stackrel{(3)}{\leq} \frac{1}{q} \cdot \sum_{i \in [1, m]} p_i \cdot H(Y_{j^{(i)}})$$

$$\stackrel{(2)}{=} \sum_{i \in [1, m]} x(i) \cdot H(Y_{j^{(i)}})$$

$$\stackrel{(4)}{\leq} \sum_{i \in [1, m]} x(i) \cdot \log N_i$$

Somit ist

$$\begin{aligned} |\mathcal{Q}(D)| &= 2^{\log(|\mathcal{Q}(D)|)} \leq 2^{\sum_{i=1}^m x(i) \cdot \log N_i} \\ &= \prod_{i=1}^m 2^{x(i) \cdot \log N_i} = \prod_{i=1}^m N_i^{x(i)} \end{aligned}$$

Die AGM-Schranke ist im folgenden Sinn optimal:

### Satz 2.10 (Optimalität der AGM-Schranke)

Sei  $Q$  eine Join-Anfrage der Form

$$Q(\bar{X}_0) \leftarrow R_1(\bar{X}_1), \dots, R_m(\bar{X}_m).$$

Für jedes  $N \in \mathbb{N}$  gibt es eine Datenbank  $D$  vom Schema  $\{R_1, \dots, R_m\}$  mit  $N_i := |R_i^D| \geq N$  für alle  $i \in \{1, \dots, m\}$  und eine fraktionale Kantenüberdeckung  $x$  von  $Q$ , so dass

$$|Q(D)| = \prod_{i=1}^m N_i^{x(i)} \quad \text{ist.}$$

Beweis:

Zum Beweis nutzen wir das aus der linearen Optimierung bekannte Dualitätsprinzip:

Zur Erinnerung:

Primales LP: (P)

$$\begin{aligned} &\text{minimiere } c^T x \\ &\text{unter der Bedingung:} \\ &Ax \geq b, \quad x \geq 0 \end{aligned}$$

$$\text{mit } c = \begin{pmatrix} c_1 \\ \vdots \\ c_m \end{pmatrix} \in \mathbb{R}^m, \quad x = \begin{pmatrix} x_1 \\ \vdots \\ x_m \end{pmatrix}$$

$$A = (a_{ji})_{\substack{j \in \{1, \dots, m\} \\ i \in \{1, \dots, n\}}}, \quad b = \begin{pmatrix} b_1 \\ \vdots \\ b_n \end{pmatrix} \in \mathbb{R}^n$$

Duales LP: (D)

$$\begin{aligned} &\text{maximiere } y^T b \\ &\text{unter der Bedingung} \\ &y^T A \leq c^T, \quad y \geq 0 \end{aligned}$$

$$\text{mit } y = \begin{pmatrix} y_1 \\ \vdots \\ y_m \end{pmatrix}$$

Für jedes  $x \in \mathbb{R}^m$  mit  $Ax \geq b, x \geq 0$   
und jedes  $y \in \mathbb{R}^n$  mit  $y^T A \leq c^T, y \geq 0$  gilt:

$$c^T x \geq y^T A x \geq y^T b.$$

Insbes. gilt für jede optimale Lösung  $x^*$  von (P)  
und jede optimale Lösung  $y^*$  von (D), dass

$$c^T x^* \geq y^{*T} b$$

Der starke Dualitätssatz besagt, dass sogar gilt:

$$c^T x^* = y^{*T} b.$$

Um Satz 2.10 zu beweisen, nutzen wir den starken  
Dualitätssatz wie folgt:

Gegeben seien eine Join-Anfrage  $Q(A_1, \dots, A_n) \leftarrow R_1(X_1), \dots, R_m(X_m)$   
und Zahlen  $N_1, \dots, N_m \in \mathbb{N}_{\geq 1}$ .

Betrachte das Primale LP  $P(Q, N_1, \dots, N_m)$ :

minimiere  $c^T x$   
unter der Bedingung  
 $Ax \geq b, x \geq 0$

$$\sum_{i=1}^n x_i \cdot \log N_i$$

f.a.  $j \in [1, m]$  ist  $\sum_{\substack{i \in [1, n]: \\ A_{ij} \in \{X_i\}}} x_i \geq 1$  und  
f.a.  $i \in [1, n]$  ist  $x_i \geq 0$

mit  $c = \begin{pmatrix} c_1 \\ \vdots \\ c_m \end{pmatrix}$  und  $c_i := \log N_i \quad \forall i \in [1, m]$ ,

$A = (a_{ji})_{\substack{j \in [1, m], \\ i \in [1, n]}}$  mit  $a_{ji} = \begin{cases} 1 & \text{falls } A_j \in \{X_i\} \\ 0 & \text{sonst} \end{cases}$

und  $b = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$ .



Das zugehörige Duale LP  $|D(Q, N_1, \dots, N_m)|$  ist dann: 53

maximiere  $y^T b$   
 unter der Bedingung  
 $y^T A \leq c^T, \quad y \geq 0$

$$\sum_{j=1}^m y_j$$

f.a.  $i \in [1, m]$  ist  $\sum_{\substack{j \in [1, m]: \\ A_{ij} \in \{\bar{x}, 1\}}} y_j \leq \log N_i$   
 und  
 f.a.  $j \in [1, m]$  ist  $y_j \geq 0$

Der starke Dualitätssatz besagt, dass für jede optimale Lösung  $x^*$  von  $P(Q, N_1, \dots, N_m)$  und jede optimale Lösung  $y^*$  von  $D(Q, N_1, \dots, N_m)$  gilt.

$$c^T x^* = y^{*T} b$$

$$\sum_{i=1}^m x_i^* \cdot \log N_i = \sum_{j=1}^m y_j^*$$

Außerdem ist laut Akh-Schranke (Satz 2.9)

$$|Q(D)| \leq \prod_{i=1}^m N_i^{x_i^*} = 2^{\sum_{i=1}^m x_i^* \cdot \log N_i} = 2^{\sum_{j=1}^m y_j^*}$$

für jede DB  $D$  mit  $|R_i^D| = N_i$  f.a.  $i \in [1, m]$ .

Wir betrachten dies hier für den Fall, dass alle Zahlen  $N_i$  Zweierpotenzen sind, d.h.  $N_i = 2^{L_i}$  mit  $L_i \in \mathbb{N}$  f.a.  $i \in [1, m]$ . Dann ist

$c_i = \log N_i = L_i \in \mathbb{N}$  f.a.  $i \in [1, m]$ , und

aus der Linearen Optimierung ist bekannt, dass das

Optimierungsproblem  $D(Q, N_1, \dots, N_m)$  eine optimale

Lösung  $y^* = \begin{pmatrix} y_1^* \\ \vdots \\ y_m^* \end{pmatrix}$  besitzt mit  $y_j^* \in \mathbb{Q}$  f.a.  $j \in [1, m]$ .

Seien  $p_1, \dots, p_m, q \in \mathbb{N}$  mit  $q \neq 0$ , s.d.  $y_j^* = \frac{p_j}{q}$  f.a.  $j \in [1, m]$ . 54

Behauptung:  $\hat{y} := \begin{pmatrix} p_1 \\ \vdots \\ p_m \end{pmatrix}$  ist eine optimale Lösung von  $D(Q, N_1^q, \dots, N_m^q)$

Beweis: Wir wissen, dass  $y^* = \begin{pmatrix} p_1/q \\ \vdots \\ p_m/q \end{pmatrix}$  eine optimale Lösung

von  $D(Q, N_1, \dots, N_m)$  ist,

maximiere  $\sum_{j=1}^m y_j$  unter der Bedingung  $y \geq 0$ ,  $y^T A \leq c^T$

mit  $c = \begin{pmatrix} c_1 \\ \vdots \\ c_m \end{pmatrix}$  und  $c_i = \log N_i$  f.a.  $i \in [1, m]$ .

Außerdem ist  $D(Q, N_1^q, \dots, N_m^q)$  das Optimierungsproblem

maximiere  $\sum_{j=1}^m y_j^q$  unter der Bedingung  $y^q \geq 0$ ,  $y^{qT} A \leq c^{qT}$

mit  $c^q = \begin{pmatrix} c_1^q \\ \vdots \\ c_m^q \end{pmatrix}$  und  $c_i^q = \log(N_i^q) = q \cdot \log N_i = q \cdot c_i$  f.a.  $i \in [1, m]$ .

Daher gilt folgendes f.a.  $y = \begin{pmatrix} y_1 \\ \vdots \\ y_m \end{pmatrix} \in \mathbb{R}^m$  und  $y^q := \begin{pmatrix} q y_1 \\ \vdots \\ q y_m \end{pmatrix}$

$y \geq 0$  und  $y^T A \leq c^T$  ( $\Rightarrow$ )  $y^q \geq 0$  und  $y^{qT} A \leq c^{qT}$

Da  $y^*$  eine optimale Lösung von  $D(Q, N_1, \dots, N_m)$  ist, ist also

$\hat{y}$  eine optimale Lösung von  $D(Q, N_1^q, \dots, N_m^q)$ .  $\square_{\text{Beh}}$

Gemäß starkem Dualitätssatz und AAK-Schranke gilt

für jede DB  $D$  mit  $|R_i^D| = N_i^q$  f.a.  $i \in [1, m]$ , dass

$$|Q(D)| \leq 2^{\sum_{j=1}^m \hat{y}_j} = 2^{\sum_{j=1}^m p_j} \text{ ist.}$$

Wir konstruieren nun eine konkrete DB  $D$ , für die gilt:

$$|R_i^D| = N_i^q \text{ f.a. } i \in [1, m] \text{ und } |Q(D)| = 2^{\sum_{j=1}^m p_j} = \prod_{j=1}^m 2^{p_j}$$

Klar: Dies beweist dann insbes die Aussage von Satz 2.10. 55

Zur Konstruktion von  $D$  gehen wir wie folgt vor:

Für jedes  $i \in [1, m]$  sei  $r_i := \text{ar}(R_i)$ , und seien  $j(i, 1), \dots, j(i, r_i) \in [1, n]$  so dass  $\bar{X}_i = A_{j(i, 1)}, \dots, A_{j(i, r_i)}$ .

D.h. das  $i$ -te Atom  $R_i(\bar{X}_i)$  im Rumpf von  $Q$  ist

von der Form  $R_i(A_{j(i, 1)}, \dots, A_{j(i, r_i)})$ .

Sei  $J(i) := \{j(i, 1), \dots, j(i, r_i)\} = \{j \in [1, n] : A_j \in \{\bar{X}_i\}\}$ .

Da  $\hat{y} = \begin{pmatrix} p_1 \\ \vdots \\ p_m \end{pmatrix}$  eine Lösung von  $D(Q, N_1^q, \dots, N_m^q)$  ist, gilt

insbes.:  $\hat{y}^T A \leq c^T$ , dh f.a.  $i \in [1, m]$  ist

$$\sum_{\substack{j \in [1, n]: \\ A_j \in \{\bar{X}_i\}}} \hat{y}_j \leq \log(N_i^q)$$

$$\sum_{k=1}^{r_i} \hat{y}_{j(i, k)} = \sum_{k=1}^{r_i} p_{j(i, k)}$$

□

Wähle die DB  $D'$  mit  $R_i^{D'} := [1, 2^{p_{j(i, 1)}}] \times \dots \times [1, 2^{p_{j(i, r_i)}}]$ .

Dann ist  $Q(D') = \left\{ t = (t_1, \dots, t_n) : \begin{array}{l} (t_{j(i, 1)}, \dots, t_{j(i, r_i)}) \in R_i^{D'} \\ \text{f.a. } i \in [1, m] \end{array} \right\}$

$$= [1, 2^{p_1}] \times \dots \times [1, 2^{p_n}],$$

und somit ist  $|Q(D')| = \prod_{j=1}^n 2^{p_j}$ .

Außerdem gilt für jedes  $i \in [1, m]$ , dass

$$|R_i^{D'}| = \prod_{k=1}^{r_i} 2^{p_{j(i,k)}} = 2^{\sum_{k=1}^{r_i} p_{j(i,k)}} \stackrel{\text{①}}{\leq} 2^{\log(N_i^q)} = N_i^q \quad \text{56}$$

Durch Hinzufügen von weiteren Tupeln erhalten wir eine DB  $D$  mit  $R_i^D \supseteq R_i^{D'}$  und  $|R_i^D| = N_i^q$  f.a.  $i \in [1, m]$ .

Es gilt:  $Q(D) \supseteq Q(D')$ , also

$$|Q(D)| \geq \prod_{j=1}^n 2^{p_j} \quad \text{— und auf Grund der}$$

AGM-Schranke (Satz 2.9) wissen wir auch, dass  $|Q(D)| \leq \prod_{j=1}^n 2^{p_j}$

ist. Dies beendet den Beweis von Satz 2.10

□

## 2.3 Verallgemeinerung der AGM-Schranke für beliebige konjunktive Anfragen

Ziel dieses Abschnitts ist, eine Verallgemeinerung der AGM-Schranke zu finden, die nicht nur für Join-Anfragen, sondern für beliebige konjunktive Anfragen gilt. Wir gehen dazu in mehreren Schritten vor, in denen wir immer allgemeinere Varianten von Anfragen betrachten.

Zunächst betrachten wir:

### Projektionen von Join-Anfragen

Definition 2.11:

(a) Eine Projektion einer Join-Anfrage ist eine konjunktive Anfrage  $Q$  der Form

$$Q(A_1, \dots, A_k) \leftarrow R_1(\bar{X}_1), \dots, R_m(\bar{X}_m),$$

für die es ein  $n \geq k$  und Variablen  $A_1, \dots, A_n$  gibt, so dass die Anfrage  $\tilde{Q}$  mit

$$\tilde{Q}(A_1, \dots, A_n) \leftarrow R_1(\bar{X}_1), \dots, R_m(\bar{X}_m)$$

eine Join-Anfrage ist.

(b) Sei  $Q(A_1, \dots, A_k) \leftarrow R_1(\bar{X}_1), \dots, R_m(\bar{X}_m)$  eine Projektion einer Join-Anfrage.

Eine fraktionale Kantenüberdeckung von  $Q$  ist eine Abbildung  $x: [m] \rightarrow \mathbb{Q}_{\geq 0}$ , s.d. f.a.

$$j \in [k] \text{ gilt: } \sum_{\substack{i \in [m] \text{ mit} \\ A_j \in \bar{X}_i}} x(i) \geq 1$$

Unter Verwendung der AGM-Schranke (Satz 2.9 & 2.10) können wir leicht folgern, dass Folgendes gilt:

Satz 2.12:

Sei  $Q$  eine Projektion einer Join-Anfrage.

(a) Für jede fraktionale Kantenüberdeckung  $x$  von  $Q$  und jede Datenbank  $D$  vom Schema  $\{R_1, \dots, R_m\}$  gilt

$$|Q(D)| \leq \prod_{i=1}^m N_i^{x(i)},$$

wobei  $N_i := |R_i^D|$  f.a.  $i \in [m]$  ist.

(b) Für jedes  $N \in \mathbb{N}$  gibt es eine  $\{R_1, \dots, R_m\}$ -DB  $D$  mit  $N_i := |R_i^D| \geq N$  f.a.  $i \in [m]$  und eine fraktionale Kantenüberdeckung  $x$  von  $Q$ , s.d.

$$|Q(D)| = \prod_{i=1}^m N_i^{x(i)}.$$

Beweis:

(a) Sei  $\tilde{Q}(A_1, \dots, A_n) \leftarrow R_1(\bar{X}_1), \dots, R_m(\bar{X}_m)$  eine zu  $Q$  gehörige Join-Anfrage (i.S.v. Def 2.11 (a)).

Für jedes  $i \in [m]$  sei  $\bar{X}'_i$  die Variablenliste, die aus  $\bar{X}_i$  entsteht, indem man alle Variablen aus  $\{A_1, \dots, A_n\}$  weglässt, sei  $r'_i$  die Länge der Liste  $\bar{X}'_i$  und sei  $R'_i$  ein neues Relationssymbol der Stelligkeit  $r'_i$ .

Wir betrachten die Anfrage  $Q'$  vom Schema  $\{R'_1, \dots, R'_m\}$  mit

$$Q'(A_1, \dots, A_k) \leftarrow R'_1(\bar{X}'_1), \dots, R'_m(\bar{X}'_m).$$

Offensichtlicherweise ist  $Q'$  eine Join-Anfrage, und für jede Abbildung  $x: [m] \rightarrow \mathbb{Q}_{\geq 0}$  gilt:

⊙  $x$  ist eine fraktionale Kantenüberdeckung von  $Q'$   $\Leftrightarrow$   $x$  ist eine fraktionale Kantenüberd. von  $Q$ .

Sei nun  $D$  eine beliebige DB vom Schema  $\{R_1, \dots, R_m\}$ .

Sei  $D'$  die DB vom Schema  $\{R'_1, \dots, R'_m\}$  mit

$$R'_i{}^{D'} = Q_i(D)$$

wobei  $Q_i$  die Anfrage  $Q_i(\bar{X}'_i) \leftarrow R_i(\bar{X}_i)$  ist, f.a.  $i \in [m]$  (d.h.  $R'_i{}^{D'}$  ist die Projektion von  $R_i^D$  auf die Komponenten, die zu  $A_1, \dots, A_k$  gehören).

Klar:  $N'_i := |R'_i{}^{D'}| \leq |R_i^D| = N_i$

Man kann sich leicht davon überzeugen, dass

$$Q(D) \subseteq Q'(D')$$

ist. Sei  $x$  eine fraktionale Kantenüberdeckung von  $Q$  (und  $Q'$ , wegen  $\textcircled{D}$ )  
 Gemäß AGM-Schranke (Satz 2.9) gilt

$$|Q'(D')| \leq \prod_{i=1}^m N_i^{x(i)}$$

$$\leq$$

$$|Q(D)|$$

$$\leq \prod_{i=1}^m N_i^{x(i)}$$

$D(a)$

(b) Wie nutzen die gleiche Notation wie im Beweis von (a). Sei  $N \in \mathbb{N}$  beliebig. Gemäß Optimalität

der AGM-Schranke (Satz 2.10) gibt es eine DB  $D'$  vom Schema  $\{R_1, \dots, R_m\}$  mit

$$N_i := |R_i^{D'}| \geq N \quad \text{f. a. } i \in [m]$$

und eine fraktionale Kantenüberdeckung  $x$  von  $Q'$  (und  $Q$ , wegen  $\textcircled{D}$ ), so dass

$$|Q'(D')| = \prod_{i=1}^m N_i^{x(i)}$$

Sei  $a$  ein beliebiges, fest gewähltes Element in  $D$

und sei  $D$  die DB vom Schema  $\{R_1, \dots, R_m\}$ , die

man aus  $D'$  erhält, indem man jedes Tupel

$t \in R_i^{D'}$  an geeigneten Stellen um weitere



Komponenten anreichert, deren Beitrag jeweils den Wert  $a$  bekommt

61

Beispiel: Falls das  $i$ -te Atom von  $Q$  die Form

$R_i(A_{z_1}, A_{k+1}, A_n)$  hat (und  $k \geq 3$ ), so

hat das  $i$ -te Atom von  $Q'$  die Form  $R_i'(A_{z_1}, A_n)$ .

Für jedes Tupel  $t' = (x, y) \in R_i'^{D'}$  enthält  $R_i^D$

dann das Tupel  $\hat{t}' := (x, a, y)$ .

Präzise: Für  $i \in [m]$  sei  $R_i(\bar{x}_i)$  von der Form

$R_i(A_{j(i,1)}, \dots, A_{j(i,r_i)})$  und sei  $M_i := \{l \in \{1, \dots, r_i\} : j(i,l) > k\}$ .

Jedem Tupel  $t' = (t'_1, \dots, t'_{r_i}) \in \text{dom}^{r_i}$  ordnen wir

das Tupel  $\hat{t}' = (\hat{t}'_1, \dots, \hat{t}'_{r_i}) \in \text{dom}^{r_i}$  zu, für das

gilt: 1)  $\hat{t}'_l = a$  f.a.  $l \in M_i$ , und

2)  $t'$  ist das Tupel, das aus  $\hat{t}'$  entsteht, indem man für jedes  $l \in M_i$  die  $l$ -te Komponente  $\hat{t}'_l$  löscht.

$D$  ist dann die DB von Schema  $\{R_1, \dots, R_m\}$  mit

$$R_i^D := \{ \hat{t}' : t' \in R_i'^{D'} \} \quad \text{f.a. } i \in [m].$$

Dann ist  $|R_i^D| = |R_i'^{D'}| = N_i$  f.a.  $i \in [m]$ .

Außerdem kann man sich leicht davon überzeugen,

dass  $Q(D) = Q'(D')$  ist.

Insgesamt gilt also:  $|Q(D)| = \prod_{i=1}^m N_i^{x(i)}$  und

$N_i = |R_i^D|$  f.a.  $i \in [m]$ .

$D(6)$

### Beispiel 2.13:

Die Anfrage

$$Q(A, B) \leftarrow E_1(A, B), E_2(B, C), E_3(C, A)$$

ist eine Projektion einer Join-Anfrage  
— nämlich der Join-Anfrage

$$\tilde{Q}(A, B, C) \leftarrow E_1(A, B), E_2(B, C), E_3(C, A).$$

Natürlich ist  $x_1: \{1, 2, 3\} \rightarrow \mathbb{Q}_{\geq 0}$  mit  
 $x_1(1) = x_1(2) = x_1(3) = \frac{1}{2}$  eine fraktionale  
Kantenüberdeckung von  $\tilde{Q}$  und von  $Q$ .

Aber gemäß Definition 2.11 ist auch  $x_2: \{1, 2, 3\} \rightarrow \mathbb{Q}_{\geq 0}$   
mit  $x_2(1) = 1$  und  $x_2(2) = x_2(3) = 0$  eine  
fraktionale Kantenüberdeckung von  $Q$   
(aber nicht von  $\tilde{Q}$ ).

Satz 2.12 liefert, dass für jede DB  $D$  gilt:

$$|Q(D)| \leq \prod_{i=1}^3 N_i^{x_2(i)} = N_1,$$

für  $N_i := |R_i^D|$  f.a.  $i \in [3]$  ist.

(... was wir bei dieser einfachen Anfrage  $Q$   
aber auch direkt, ohne Nutzung der A&K-Schranke,  
sehen können).

## Konjunktive Anfragen ohne Konstanten

### Definition 2.14

Sei  $Q(\bar{x}_0) \leftarrow R_1(\bar{x}_1), \dots, R_m(\bar{x}_m)$  eine konjunktive Anfrage, in der keine Konstanten vorkommen (d.h.  $\text{Cons}(Q) = \emptyset$ ), sei  $\{A_1, \dots, A_n\} = \text{vars}(Q)$  und  $\{\bar{x}_0\} = \{A_1, \dots, A_k\}$ .

Eine fraktionale Kantenüberdeckung von  $Q$  ist eine Abbildung  $x: [m] \rightarrow \mathbb{Q}_{\geq 0}$ , s.d. f.a.  $j \in [k]$  gilt:

$$\sum_{\substack{i \in [m] \text{ mit} \\ A_j \in \{\bar{x}_i\}}} x(i) \geq 1$$

### Satz 2.15 (Gottlob, Lee, Valiant, 2009)

Sei  $Q$  eine konjunktive Anfrage mit  $\text{cons}(Q) = \emptyset$  der Form  $Q(\bar{x}_0) \leftarrow R_1(\bar{x}_1), \dots, R_m(\bar{x}_m)$ .

(a) Für jede fraktionale Kantenüberdeckung  $x$  von  $Q$  und jede Datenbank  $\mathcal{D}$  und alle  $N_1, \dots, N_m \in \mathbb{N}$  mit  $|R_i^{\mathcal{D}}| \leq N_i$  f.a.  $i \in [m]$  gilt:

$$|Q(\mathcal{D})| \leq \prod_{i=1}^m N_i^{x(i)}$$

(b) Für jedes  $N \in \mathbb{N}$  gibt es natürliche Zahlen  $N_1, \dots, N_m \geq N$ , eine Datenbank  $\mathcal{D}$  mit  $|R_i^{\mathcal{D}}| \leq \sum_{\substack{j \in [m] \text{ mit} \\ R_j = R_i}} N_j$  f.a.  $i \in [m]$  und eine fraktionale Kantenüberdeckung  $x$  von  $Q$ , s.d.

$$|Q(\mathcal{D})| \geq \prod_{i=1}^m N_i^{x(i)}$$

Beweis:

(I) Wir zeigen die Aussage zunächst für selbst-join-freie Anfragen, d.h. für Anfragen, bei denen die Relationssymbole  $R_1, \dots, R_m$  paarweise verschieden sind (d.h.  $R_i \neq R_j$  f.a.  $i, j \in [m]$  mit  $i \neq j$ ).

Für jedes  $i \in [0, m]$  sei  $\bar{X}_i^{\text{versch}}$  die Variablenliste, die man aus  $\bar{X}_i$  erhält, indem man Mehrfach-Vorkommen derselben Variablen entfernt. ("versch" steht für "verschieden").

Insbes. ist  $\{\bar{X}_i^{\text{versch}}\} = \{\bar{X}_i\}$ .

Sei  $k_i := |\{\bar{X}_i\}| = |\{\bar{X}_i^{\text{versch}}\}|$  und sei  $R_i^{\text{versch}}$  ein neues Relationssymbol der Stelligkeit  $k_i$ :

Betrachte die Anfrage

$$Q^{\text{versch}}(\bar{X}_0^{\text{versch}}) \leftarrow R_1^{\text{versch}}(\bar{X}_1^{\text{versch}}), \dots, R_m^{\text{versch}}(\bar{X}_m^{\text{versch}}).$$

Diese Anfrage ist eine Projektion einer Join-Anfrage. (i.S.v. Def. 2.11). Außerdem gilt für jede Abbildung

$x: [m] \rightarrow \mathbb{Q}_{\geq 0}$ :

$\Leftrightarrow$ :  $x$  ist eine frakt. Kantenüberdeckung von  $Q^{\text{versch}}$   $\Leftrightarrow$   $x$  ist eine frakt. Kantenüberdeckung von  $Q$

Betrachte eine beliebige DB  $D$  vom Schema  $\{R_1, \dots, R_m\}$ .

Sei  $D'$  die DB vom Schema  $\{R_1^{\text{versch}}, \dots, R_m^{\text{versch}}\}$  mit

$$(R_i^{\text{versch}})^{D'} := Q_i(D) \text{ mit } Q_i(\bar{X}_i^{\text{versch}}) \leftarrow R_i(\bar{X}_i).$$

Dann ist  $N_i^{D'} := |(R_i^{\text{versch}})^{D'}| \leq |R_i^D|$ .

und es gilt:  $|Q(D)| = |Q^{\text{versch}}(D')|$  (\*) 65

Gemäß Satz 2.12(a) gilt für jede fraktionale  
Kantenüberdeckung  $x$  von  $Q^{\text{versch}}$  (und  $Q$ , wegen  $\textcircled{D}$ ),

$$\text{dass } |Q^{\text{versch}}(D')| \leq \prod_{i=1}^m N_i^{x(i)}.$$

Also gilt:

$$|Q(D)| \stackrel{(*)}{=} |Q^{\text{versch}}(D')| \leq \prod_{i=1}^m N_i^{x(i)} \leq \prod_{i=1}^m |R_i^D|^{x(i)} \leq \prod_{i=1}^m N_i^{x(i)}$$

f.a.  $N_1, \dots, N_m \in \mathbb{N}$  mit  $|R_i^D| \leq N_i$  f.a.  $i \in [m]$ .

Also gilt Aussage (a) für selbst-join-freie Anfragen  $Q$ .

Zum Beweis von Aussage (b) für selbst-join-freie Anfragen  $Q$   
betrachte ein beliebiges  $N \in \mathbb{N}$ .

Gemäß Satz 2.12 (b) gibt es eine  $\{R_1^{\text{versch}}, \dots, R_m^{\text{versch}}\}$ -DB  
 $\hat{D}$  mit  $\hat{N}_i := |(R_i^{\text{versch}})^{\hat{D}}| \geq N$  f.a.  $i \in [m]$  und eine  
fraktionale Kantenüberdeckung  $x$  von  $Q^{\text{versch}}$  (und  $Q$  wg.  $\textcircled{D}$ ),

$$\text{s.d. } |Q^{\text{versch}}(\hat{D})| = \prod_{i=1}^m \hat{N}_i^{x(i)}. \quad (**)$$

Man kann leicht eine DB  $D$  vom Schema  $\{R_1, \dots, R_m\}$   
konstruieren, s.d.  $D' = \hat{D}$  und  $|R_i^D| = |(R_i^{\text{versch}})^{D'}|$

f.a.  $i \in [m]$  ist (Details: Übung). Insgesamt gilt:

$$|Q(D)| \stackrel{(*)}{=} |Q^{\text{versch}}(D')| = \prod_{i=1}^m \hat{N}_i^{x(i)} = \prod_{i=1}^m |R_i^D|^{x(i)}$$

Dies beweist Aussage (b) für selbst-join-freie Anfragen  $Q$ .

66

(II) Wir zeigen nun, dass die Aussage auch für Anfragen  $Q$  gilt, die nicht selbst-join-frei sind. Wir ordnen dazu der Anfrage

$$Q(\bar{x}_0) \leftarrow R_1(\bar{x}_1), \dots, R_m(\bar{x}_m)$$

die selbst-join-freie Anfrage

$$Q'(\bar{x}_0) \leftarrow R'_1(\bar{x}_1), \dots, R'_m(\bar{x}_m)$$

vom Schema  $\{R'_1, \dots, R'_m\}$  zu, wobei  $R'_1, \dots, R'_m$  paarweise verschiedene Relationssymbole mit  $ar(R'_i) = ar(R_i)$  f.a.  $i \in [m]$  sind.

Gemäß (I) wissen wir, dass die Aussagen (a) und (b) für die Anfrage  $Q'$  gelten. Außerdem gilt für jede Abbildung  $x: [m] \rightarrow \mathbb{Q}_{\geq 0}$ :

$\textcircled{\triangleright}$   $x$  ist eine fraktionale Kantensbedeckung von  $Q$   $\Leftrightarrow$   $x$  ist eine fraktionale Kantensbedeckung von  $Q'$ .

Zum Beweis von (a) ordnen wir jeder  $\{R_1, \dots, R_m\}$ -DB  $D$  die  $\{R'_1, \dots, R'_m\}$ -DB  $D'$  mit  $R'_i D' := R_i D$  f.a.  $i \in [m]$  zu. Klar:  $Q(D) = Q'(D')$ .

Damit erhalten wir, dass Aussage (a) auch für  $Q$  gilt.

Zum Beweis von Aussage (b) für  $Q$  betrachte ein beliebiges  $N \in \mathbb{N}$ . Gemäß (I)(b) gibt es eine

$\{R'_1, \dots, R'_m\}$ -DB  $D'$  mit  $N_i := |R'_i D'| \geq N$  f.a.  $i \in [m]$  und eine fraktionale Kantensbedeckung  $x$  von  $Q'$  (und  $Q$ , wegen  $\textcircled{\triangleright}$ ), s.d.  $|Q'(D')| \geq \prod_{i=1}^m N_i^{x(i)}$  (Begründung: Übung!)

Sei  $\tilde{D}$  die DB vom Schema  $\{R_1, \dots, R_m\}$  mit

$$R_i^{\tilde{D}} := \bigcup_{\substack{j \in [m] \text{ mit} \\ R_j = R_i}} R_j^{D'}$$

Man kann leicht nachprüfen, dass

$$Q(\tilde{D}) \supseteq Q'(D').$$

Also gilt:

$$|Q(\tilde{D})| \geq |Q'(D')| \geq \prod_{i=1}^m N_i^{x(i)}$$

Anßerdem gilt  $\forall a. i \in [m]$ , dass

$$|R_i^{\tilde{D}}| \leq \sum_{\substack{j \in [m]: \\ R_j = R_i}} |R_j^{D'}| = \sum_{\substack{j \in [m] \\ R_j = R_i}} N_j$$

Dies zeigt, dass Aussage (b) für  $Q$  gilt

□