

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

50
wobei $Y_{j^{(i)}} = (Y_j)_{j \in J_k^{(i)}} = (Y_j)_{A_j \in \{X_i\}}$ ist

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

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

□ Satz 2.9

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 $\sum_{\substack{i \in [1, n]: \\ A_{ij} \in \{T_i\}}} x_i \geq 1$ 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 $a_{ji} = \begin{cases} 1 & \text{falls } A_j \in \{T_i\} \\ 0 & \text{sonst} \end{cases}$

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 A&K-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_{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

□