

分散データベースシステムにおける 一般化準結合を用いた質問処理

京都大学 工学部 吉川 正俊
上林 弥彦
矢島 脩三

1. まえがき

分散型データベースにおける質問処理では、全質問処理コストのうちデータ転送コストが支配的であるため、質問処理手順を構築する上では、各地点間のデータ転送量を可能な限り削減することが重要である⁽¹⁾。異なる地点に存在する関係を結合する場合、データ転送量削減の上で有効な手法として準結合 (Semi-join) 操作が知られている。

質問処理の最適化を考える上で、等結合質問は重要なクラスである。等結合質問は木型質問と巡回型質問に二分され、このうち木型質問は一般に関係の数を n とすると $2(n-1)$ 回の準結合操作を実行することによりすべての関係を解の状態にできることが知られている⁽²⁾⁽³⁾。それに対し巡回型質問は、一般に準結合操作のみを用いて解くことは不可能であり、またたとえ部分的に質問を解く場合でも、必要となる準結合操

作の回数が非常に多くなる⁽²⁾⁽³⁾などの問題があり、現在のところ巡回型質問を効率よく解くための手法は見い出されていない。

このため本稿では準結合操作を一般化準結合操作を一般化した一般化準結合操作を提案する。一般化準結合操作は、2つの関係間の結合属性値に加えて、必要な属性の値をと同時に転送を行なう操作である。質問グラフの全域木を考えることにより、各関係間でどの属性の値の転送が必要かがわかる。またこの全域木に沿って一般化準結合操作を実行することにより、任意の巡回型質問を解くことができることを示す。

2. 諸定義

属性 A_1, A_2, \dots, A_m からなる関係 R とは各属性 $A_i (i=1, 2, \dots, m)$ の定義域の直積の部分集合である。 R の各要素を組と言う。関係 R の属性集合を関係スキーマと言い $R[A_1, A_2, \dots, A_m]$ あるいは R でそれを表わす。関係スキーマの集まりをデータベーススキーマと言い $D[R_1, R_2, \dots, R_n]$ で表わす。データベース状態 D とは、各関係の直積とする。すなわち、

$$D = R_1 \times R_2 \times \dots \times R_n$$

$A_{i \in R}, A_{j \in L}$ をそれぞれ R_i, R_j の属性とするとき $R_i.A_{i \in R} = R_j.A_{j \in L}$ の形をした条件節の論理積結合を等結合条件式と言う。またすべての条件節 $R_i.A_{i \in R} = R_j.A_{j \in L}$ において属性 $A_{i \in R}$ と $A_{j \in L}$ が等

しいような等結合条件式を半自然結合条件式と言う。さらに、データベーススキーマ \mathcal{D} において2つの属性 A_{iR} と A_{jL} が等しいとき、またそのときに限り条件節 $R_i.A_{iR} = R_j.A_{jL}$ を持つような等結合条件式を(\mathcal{D} における)自然結合条件式と言う。条件式 ξ に対応する質問とはデータベース状態 \mathcal{D} をデータベース状態 $\{d \in \mathcal{D} \mid \xi(d) \text{ is true}\}$ に写す関数であるとする。以後、本稿では質問も条件式と同様に ξ で表わすものとする。(したがって $\{d \in \mathcal{D} \mid \xi(d) \text{ is true}\}$ も単に $\xi(\mathcal{D})$ で表わす。)等結合条件式に対応する質問を等結合質問と言い(半)自然結合条件式に対応する質問を(半)自然結合質問と言う。定義から明らかなように自然結合質問は与えられたデータベーススキーマに対して唯一つ存在する。2つの質問は任意のデータベース状態に対する答が等しいとき等価であると言う。任意の等結合質問は属性名変更することにより自然結合質問に変換することができる⁽³⁾。

本稿では質問として自然結合質問を仮定する。

半自然結合質問 ξ に対する質問グラフ $G_\xi(V_\xi, E_\xi, L_\xi)$ はラベル付きの無向グラフである。ここで節点集合 V_ξ は ξ において参照される関係の集合を表わし、質問中に $R_i.A = R_j.A$ なる条件節が存在するときまたそのときに限り節点 R_i と R_j の間をラベル A を持つ枝で結ぶ。

質問グラフが木で表わされる半自然結合質問あるいはそれと等価な質問を木型質問と言い，木型質問以外の半自然結合質問を巡回型質問と言う。

[例1] データベーススキーマ $D[R_1[A, B, C], R_2[A, C, D], R_3[A, B, E]]$ における以下の3つの質問 Q_1, Q_2, Q_3 の質問グラフをそれぞれ図1の(a) ~ (c)に示す。

$$Q_1: (R_1.A = R_2.A) \wedge (R_2.B = R_3.B) \wedge (R_3.C = R_1.C)$$

$$Q_2: (R_1.A = R_2.A) \wedge (R_2.B = R_3.B) \wedge (R_3.A = R_1.A)$$

$$Q_3: (R_1.A = R_2.A) \wedge (R_2.A = R_3.A) \wedge (R_2.B = R_3.B)$$

質問 Q_2 はその質問グラフに閉路を持つが，質問 Q_3 と等価であるため木型質問である。それに対し質問 Q_1 は巡回型質問である。□

R_1 と R_2 の自然結合は $R_1 \bowtie R_2$ で表わす。また R_1, R_2, \dots, R_n の n 個の関係の自然結合は簡単に $\bowtie R_i$ で表わす。

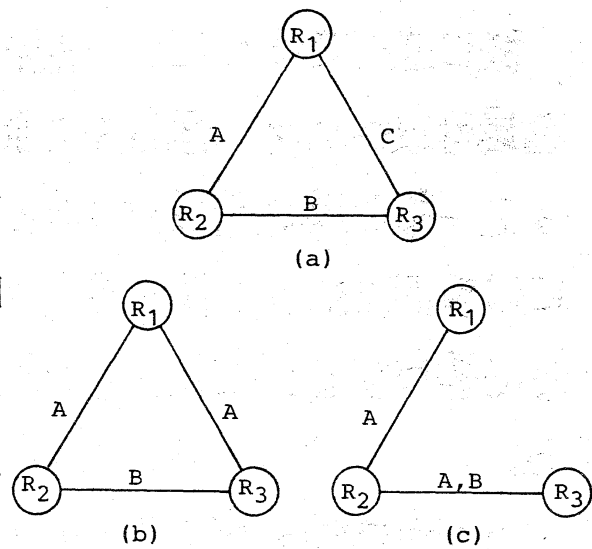


図1. 質問グラフ

R_2 による R_1 の準結合は $R_1 \ltimes R_2$ で表わされ，次式で定義される：
 $R_1 \ltimes R_2 \triangleq (R_1 \bowtie R_2)[R_1]$

X を $R_1 \cap R_2 \subseteq X \subseteq R_2$ を満足する属性集合とするとき $R_1 \ltimes (R_2[X])$ を X における R_2 による R_1 の一般化準結合と言い，それを $R_1 \ltimes^* R_2$

で表わす。 $R_1 \cap R_2 = X$ のときは $R_1 \boxtimes R_2 = R_1 \times R_2$ が成立する。
一般化準結合を実行することにより関係スキーマ自体が動的
に変化する。

本稿では分散型データベースを仮定し、各関係は互いに異
なる一つの地点に存在するものとする。

3. 一般化準結合を用いた質問処理

3.1 質問全域木

本節では、任意の巡回型質問を一般化準結合を用いて解く
方法について述べる。

既に述べたように、一般化準結合操作は2つの関係間で結
合属性以外の他の属性をと転送を行なう操作である。したが
って一般化準結合操作を用いた質問処理手順を構築するため
には質問に応じて関係間でどの属性値を転送するかを決定す
る必要がある。次に、まずそのための写像として $attr$ を与え
る。

巡回型質問 q が与えられたとき質問グラフ G_q の根付き全域
木を質問全域木と言い、それを T_q で表わす。 T_q は木であるた
め、節点間の親子関係が定義できる。 T_q において節点 R_i とそ
の親の節点を結ぶ枝を a_i で表わし、節点 R_i とその子孫から成
る節点の集合を α_i で表わすものとする。一般性を失うことな
く T_q の根を R_1 とする。 T_q の枝の集合を E_q とするとき、 T_q の節

及び枝から属性集合への写像 $attr$ を図2に示す手続きによって定義する。ただし $label(e)$ は枝 $e (e \in E_g)$ に付いている属性ラベルを表わすものとする。 S を質問中に現われるすべての結合属性の集合とすると、図2の手続きは $O(|E_g| |V_g| |S|)$ の時間ですべての枝及び節の $attr$ の値を計算する。 $attr$ の定義より次の命題が成立する。

[命題1] 節点 R_i の親を R_j とするとき

$$attr(a_i) = attr(R_i) \cap attr(R_j)$$

do for each $e \in E'_q$;

$attr(e) = \{label(e)\}$;

end;

do for each $R_i \in V_q$;

$attr(R_i) = \underline{R}_i$;

end;

do for each $e = (R_j, R_k) \in E_q - E'_q$;

do for each $e' \in T_q$ におけるパス (R_j, R_k) に含まれる枝の集合;

$attr(e') = attr(e') \cup \{label(e)\}$;

end;

do for each $R_i \in T_q$ におけるパス (R_j, R_k) に含まれる節の集合;

$attr(R_i) = attr(R_i) \cup \{label(e)\}$;

end;

end;

図2 $attr$ を計算する手続き

[例2] 図3に質問全域木の例を示す。各節及び枝に写像attrを適用した結果は以下のようなになる。

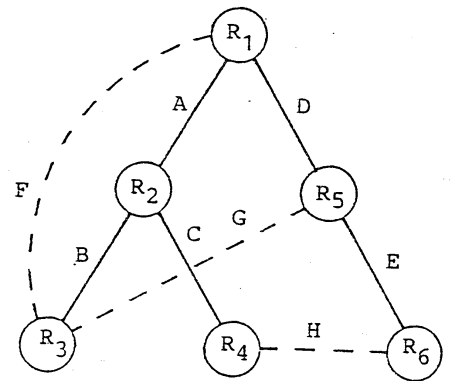


図3 質問全域木

$$\text{attr}(R_1) = \{A, D, F, G, H\}$$

$$\text{attr}(R_2) = \{A, B, C, F, G, H\}$$

$$\text{attr}(R_3) = \{B, F, G\}$$

$$\text{attr}(R_4) = \{C, H\} \quad \text{attr}(R_5) = \{D, E, G, H\}$$

$$\text{attr}(R_6) = \{E, H\}$$

$$\text{attr}(a_2) = \{A, F, G, H\} \quad \text{attr}(a_3) = \{B, F, G\}$$

$$\text{attr}(a_4) = \{C, H\} \quad \text{attr}(a_5) = \{D, G, H\} \quad \text{attr}(a_6) = \{E, H\} \quad \square$$

$(\bigotimes_{j=1}^n R_j) [\text{attr}(R_i)]$ を R_i の既約という。既約な関係は与えられた自然結合条件式を真にするような最小の状態になっている。次に一般化準結合を用いて各関係を既約にするための手続きを示す。

3.2 質問処理手続きUP

$T_{\mathcal{Q}}$ において節 R_i が m 個 ($m \geq 0$) の子 R_{i1}, \dots, R_{im} を持つとする。このとき、 R_i に対する手続き $UP(\mathcal{Q}, R_i)$ を以下のように再帰的に定義する。

(i) $m = 0$ のとき $UP(\mathcal{Q}, R_i) = \langle \quad \rangle$

(ii) $m \neq 0$ のとき

$$UP(\mathcal{Q}, R_i) = \langle UP(\mathcal{Q}, R_{i1}), \dots, UP(\mathcal{Q}, R_{im}), R_i \otimes_{\text{attr}(a_{i1})} R_{i1}, \dots, R_i \otimes_{\text{attr}(a_{im})} R_{im} \rangle$$

ただし、 $\langle \rangle$ 内の手続きは左から順に処理を行なうものとする。直感的には、 $UP(\xi, R_i)$ は T_ξ の葉から節 R_i に向けて順に一般化準結合操作を繰り返す手続きである。手続き $UP(\xi, R_i)$ を単に UP で表わす。

\hat{R}_i に属するすべての関係の自然結合を $\bowtie \hat{R}_i$ で表わすものとする。このとき次の補題が成立する。

[補題1] 手続き $UP(\xi, R_i)$ の結果、関係 R_i は $\bowtie \hat{R}_i [Attr(R_i)]$ となる。

(証明) 略 □

補題1よりただちに次の定理が成立する。

[定理1] 手続き UP の結果、関係 R_i は既約となる。□

3.3 全域木の根以外の関係の既約化

前項で述べた手続き UP により、全域木の根に相当する関係 R_i は既約となることがわかったが、 R_i 以外の関係は一般にまだ既約とは言えない。本稿では R_i 以外の関係を既約にするための手続きについて考察する。

[定理2] T_ξ において関係 R_i の状態が $\bowtie \hat{R}_i [Attr(R_i)]$ であり、 $attr(a_i) \leq attr(R_j)$ を満足する既約な関係 R_j が存在するならば、一般化準結合操作 $R_i \bowtie_{attr(a_i)} R_j$ により、 R_i もまた既約になる。

(証明) 略 □

T_q において R_i が R_i の親の場合は、命題1より $attr(a_i) \leq attr(R_i)$ が成立する。したがって定理2より、手続きUPを適用後以下で定義される手続きDOWNを実行することにより、すべての関係を既約にできることがわかる。手続きDOWNとUPと同様に以下に再帰的な定義を与える。

(i) $m = 0$ のとき $DOWN(\underline{g}, R_i) = \langle \quad \rangle$

(ii) $m \neq 0$ のとき $DOWN(\underline{g}, R_i) = \langle R_{i1} \boxtimes_{attr(a_{i1})} R_i, \dots, R_{im} \boxtimes_{attr(a_{im})} R_i, DOWN(\underline{g}, R_{i1}), \dots, DOWN(\underline{g}, R_{im}) \rangle$

UPと同様に手続きDOWN(\underline{g}, R_i)をDOWNで表わす。

[例3] 図4に質問グラフと質問全域木を示す。この例では $attr(R_1) = \{A, B, D\}$, $attr(R_4) = \{A, B\}$ である。したがって定理2より、手続きUPが終了した時点で一般化準結合操作 $R_4 \boxtimes_{AB} R_1$ を実行することにより、関係 R_4 は既約となる。すなわち、この例の場合、手続きUP終了後、各関係を既約にするため一般化準結合操作の系列としては

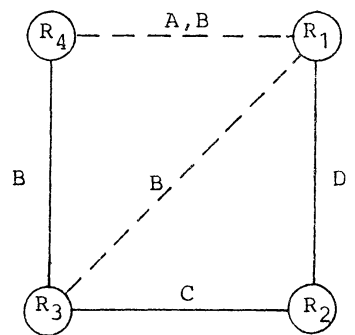


図4. 関係の既約化

$$R_2 \boxtimes_{ABD} R_1 \Rightarrow R_3 \boxtimes_{ABC} R_2 \Rightarrow R_4 \boxtimes_{AB} R_3$$

及び

$$R_4 \boxtimes_{AB} R_1, R_2 \boxtimes_{ABD} R_1 \Rightarrow R_3 \boxtimes_{ABC} R_2$$

の2通りが考えられる。□

5. あとがき

本稿では、省略したが、通常の準結合に比べて余分に転送が必要な属性値をデータ圧縮することにより一般化準結合を用いた巡回型質問の処理は通常の準結合操作と同等の転送データ量で実行が可能であると考ええる。

謝辞 御討論頂いた矢島研諸氏に深謝する。なお、本研究は一部文部省科学研究費による。

参考文献

- (1) Rothnie, J.B., Jr. and Goodman, N. "A Survey of Research and Development in Distributed Database Management," Proceedings of International Conference on VLDB, pp.48-62, Oct. 1977.
- (2) Bernstein, P.A. and Chiu, D.M. "Using Semi-Joins to Solve Relational Queries," JACM, Vol.28, No.1, pp.25-40, Jan. 1981.
- (3) Bernstein, P.A. and Goodman, N. "The Theory of Semi-Joins," CCA Rep., No.CCA-79-27, Nov.15, 1979.