

1.はじめに

演繹データベースシステムでは単一化可能パタンの検索は要求パターン・対象パターン共に複数個存在する場合、組合せが増大し処理の隘路となりやすい。本稿ではこの操作を単一化を伴う関係代数の結合演算とみなし、ハッシュとソートによって実行するアルゴリズムについて述べる。また、他方式との比較を行いその最適化アルゴリズムについても述べる。

2.単一化を伴う結合演算のアルゴリズム

以下では原子論理式を属性値に持つ関係A, B間の単一化を伴う結合演算を例にとる。簡単のために各関係はハッシング等によって述語, アリティー一定のタプル集合に予め分割されていると仮定する (Fig.1)。

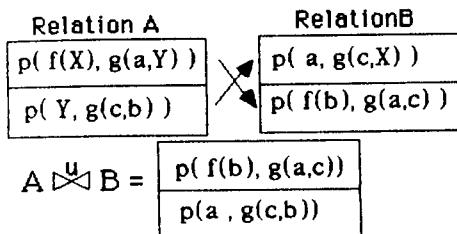


Fig.1 an example for join with unification

まずアルゴリズムの戦略について説明する。

与えられた述語のアリティーに応じてFig.2のように木構造を設定し原子論理式からなるタプル $t_1, t_2$ をこれにあわせて展開する。

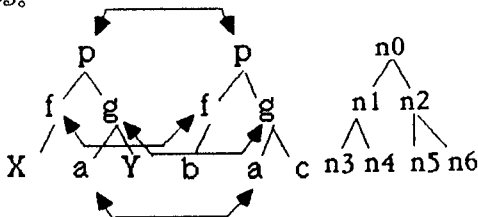


Fig.2 an example for unifiable atomic formulae  
 $\longleftrightarrow$  : equal symbol

このとき二つのタプル $t_1, t_2$ が単一化可能であればハッシュ木の各ノード $n_i$ について次の条件が成立する。

「 $n_i$ に対応するタプル $t_1, t_2$ のシンボルを $s_1, s_2$ とすると $s_1, s_2$ がともに変数でなければこれらは一致しなければならない

い (Fig.2)。」

アルゴリズムは以下の通り (Fig.3)。

STEP1. ハッシング

関係A, Bの各タプル $t$ についてハッシュ木のノード $n_0, n_1, \dots$ に応じたシンボル $s_0, s_1, \dots$ を取り、ベクタ形式のハッシュ値 (以下ハッシュベクトル)  $\{h(s_0), h(s_1), \dots\}$ を求める。

(但し、 $h(\text{変数シンボル}) = 0, h(\text{その他}) = 1 \sim N$ )

このハッシュベクトルに応じて関係A, Bをバケット集合 $\{a_i\}, \{b_j\}$  ( $i, j = 0, 1, \dots$ ) にクラスタリングする。

STEP2. ソーティング

任意のバケット組合せ $a_i, b_j$ についてハッシュベクタ  $H(a_i) = (p_1, p_2, \dots), H(b_j) = (q_1, q_2, \dots)$ とする。

このときもし $p_k \cdot q_k \neq 0$ , かつ $p_k \neq q_k$ なる $k$ が存在すれば、バケット $a_i, b_j$ の任意のタプル組合せについてノード $n_k$ に対応する非変数シンボルは一致せず故に単一化可能なタプル組合せは存在しない。そうでなければ $p_k = q_k \neq 0$ なる $k$ に対応するハッシュ木上のノード列をキーにしてバケット $a_i, b_j$ のタプルをマージソートする。その後、ソート列中のキーが等しいA, Bのタプル同士を総当たりで単一化試行。

$$\begin{aligned} p(f(X), g(a, Y)) &\in a_i \\ p(f(b), g(a, c)) &\in b_j \end{aligned}$$

$$\begin{aligned} H(a_i) &= [h(p), h(f), h(g), 0, 0, h(a), 0] \\ H(b_j) &= [h(p), h(f), h(g), h(b), 0, h(a), h(c)] \end{aligned}$$

$$\begin{matrix} 0 & 0 & 0 & \times & \times & 0 & \times \\ \text{node sequence } [n_0, n_1, n_2, n_5] & \text{makes key} \end{matrix}$$

Fig.3 an example for unifiable bucket

3. 他方式との比較

単一化を伴う結合演算のアルゴリズムはICOTからも提案されている。この方法 (以下ICOT法) では原子論理式をbreadth-firstに展開し変数>定数の全順序をつけてマルチキーソートを行う。そのため優先順位の高いキーに変数シンボルがある原子論理式はソートストリームの先頭に現れそこからストリームの最後まで単一化試行を必要とし、単一化負荷の増大を招く。

これに対しhash&sort法は予め構造情報によるクラスタリングを行うことによって変数になりやすいノードがキーになることを回避し、また、「タプルを大量に発生する原因になる、定数になりやすいノード」をキーにして絞り込む。

関係A,B間の単一化を伴う自然結合をhash&sort法, ICOT法で実行した (Table.1)。サンプルS1からS3に行くに従いbreadth-firstで優先順位の高いノードの変数確率が高くなっている。単一化負荷, ソート負荷は各々  $|A| \cdot |B|$ ,  $|A| + |B|$  に対する比率である。ICOT法が変数の出現確率に大きく影響されるのに対しhash&sort法は常により高い絞り込み効果を達成できることがわかる。一方, hash&sort法のソート負荷はICOT法に対し2~6倍となる。そのため単一化負荷, ソート負荷間のトレードオフを考慮し, タプルのパタン分布に応じたハッシュ木の最適化が必要となる。

Table 1. results of join with unification.  
value = <ICOT> / <hash&sort>  
( ) is default success ratio.

sample	load of unification	success ratio	load of sorting	$ A  \cdot  B $
S1	0.3 / 0.12	0.3 / 0.64 (0.08)	1.0 / 5.6	$10^4$
S2	0.63 / 0.42	0.61 / 0.82 (0.37)	1.0 / 5.3	$10^3$
S3	1.0 / 0.27	0.22 / 0.81 (0.22)	1.0 / 1.9	$10^3$

4. hash & sort法の最適化アルゴリズム

hash&sort法では変数になる確率が低いノードを使ってクラスタリングをしても単一化負荷はさほど軽減されずソート負荷の増大によって却って全体のコストが増大する場合がある。そうしたノード列はICOT法のように全順序をつけ, hash&sort法によるキー (以下主キー) の後ろに補助キーとして追加すれば良い (Fig.4)。ソート列中主キーが一致し, 補助キーがICOT法同様全順序を保つタプル組み合わせのみ単一化を試行する。

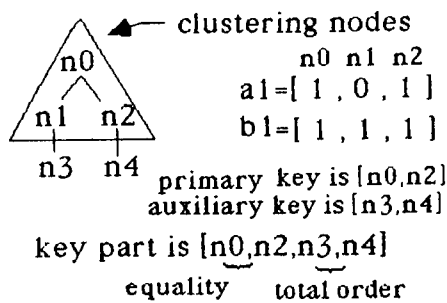


Fig.4 an optimized key

従って関係A,Bのタプルのパタン分布に応じたハッシュ木の設定, つまり変数/非変数によってクラスタリングを行うノード (以下クラスタリングノード) の決定, および上述した補助キーにするノードの決定が必要となる。

以下はそのノード列生成アルゴリズムである。

今, ノード列  $n_{i_1}, \dots, n_{i_k}$  をクラスタリングノードとする

バケット  $a_i, b_j$  を取る。  $a_i, b_j$  をさらにノード  $n_{i_k}$  の変数/非変数に応じて  $a_{i0}, a_{i1}, b_{j0}, b_{j1}$  に分割するとこれによる単一化負荷の増分  $\Delta U$ , ソート負荷の増分  $\Delta S$  は以下のようなになる (Fig.5)。

$$\Delta S = (|a_{i0}| + |b_{j0}|) + (|a_{i0}| + |b_{j1}|) + (|a_{i1}| + |b_{j0}|) + (|a_{i1}| + |b_{j1}|) - (|a_i| + |b_j|) = |a_i| + |b_j|$$

$\Delta U$  = ノード  $n_k$  をキーにすることで省かれる単一化試行数 = Fig.5 の 部にあたる負の値

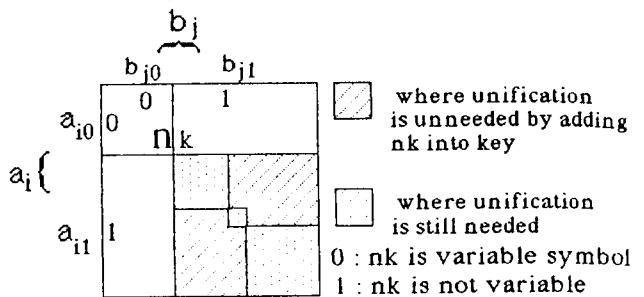


Fig.5 an example for cost-increment

故にコスト増分  $\Delta C = K_s \cdot \Delta S + K_u \cdot \Delta U$

( $K_s$  : 単位ソート負荷,  $K_u$  : 単位単一化負荷)

とすると  $\Delta C < 0$  ならノード  $n_{i_k}$  をクラスタリングノードに追加する。  $\Delta C > 0$  ならまだクラスタリングノードに含まれていないノードの内, 変数になる確率が低いものから優先順位を高くして補助キーとする。関係A,Bの直積空間をこの方法で述語シンボルから始めて再帰的にブロックに分割し, 各ブロック毎にその条件を満たすA,Bのバケット集合を取る。そしてそのブロックの指定ノードをFig.4に従ったキーとしてソートを行う (Fig.6)。

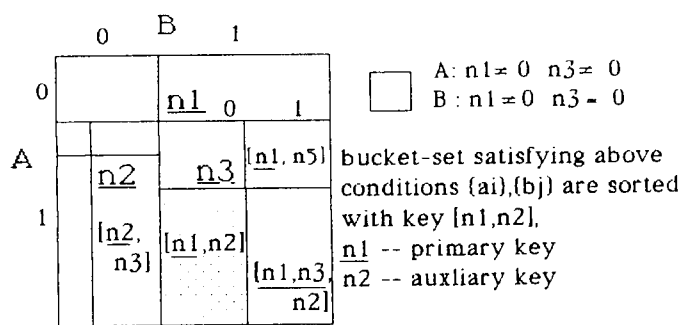


Fig.6 an example for optimize: buckets

5. おわりに

本稿では単一化を伴う結合演算をハッシュとソートによって実行するアルゴリズムについて述べ, 他方式と比較しその特性に基づいた最適化法を提案した。我々は, 多くのアプリケーションではクラスタリングノードは少数で充分と考えている。今後はhash&sort法の最適化アルゴリズムの有効性を確認していく予定である。

(参考文献)

大森 他 "推論機能と関係データベースの融合方式"

電子通信学会技術研究報告AI-86-21