

## 広域検索のための関係代数の拡張

大塚真吾 宮崎収兄

千葉工業大学情報工学科

広域ネットワークでは不特定多数のデータベースが散在している。このような環境での問合せでは問合せたい内容によって対象データが変化するため、あらかじめスキーマの統合を行うことは困難である。我々はネットワーク上にサイトが多数存在する環境で全体を一つの大きなデータベースと考え、各データベースはその一部のデータを格納した不完全なデータベースと考えることにより広域検索を行う方法を提案した。本稿では、この方法を実現するために条件付タブル、マージなどの概念を導入し関係代数を拡張して広域データベース検索処理の具体化を検討する。

## Extending Relational Algebra for Global Query Processing

Shingo Otsuka and Nobuyoshi Miyazaki

(otsuka, miyazaki)@mz.cs.it-chiba.ac.jp

Dept. of CS, Chiba Institute of Technology

2-17-1 Tsudanuma Narashino Chiba 275 Japan

It is difficult to process even simple global queries if many independent databases are scattered in a large network like the Internet, because there are various mismatches between databases. A collection of databases in a network can be regarded as a very large virtual database, where individual databases contain partial incomplete information. We proposed a method to process global queries based on the concept of incomplete database. This paper discusses extended relational algebra based on extended tuple concepts and introduces extended operations such as "merge" that are useful to realize global query processing.

## 1 はじめに

広域ネットワークでの検索はインターネット上のWWWなどの各種情報ツールによって扱われている。検索ツールの背後にデータベース管理システムを用いることも行われているが、これらはデータベースの条件検索機能を直接広域検索に用いていない。WWWのようなリンクによるアクセスやサーチエンジンによるホームページの内容のキーワード検索だけでなく、データベースに対する質問処理が実現できればネットワークを通じた情報の高度利用が可能となる。従来のマルチデータベースではあらかじめ定められた少数のデータベースの統合的扱いが課題であり、広域ネットワークのように不特定多数のデータベースが散在する環境での統合的な検索はあまり研究されていない。このような環境での問合せでは問合せたい内容によって対象データが変化するため、あらかじめスキーマの統合を行うことは困難である。このため広域ネットワークでのデータベースの質問処理を可能とするための新しい枠組みが必要となっている。我々はネットワーク上に情報提供サイトが多数存在する環境で全体を一つの大きなデータベースと考え、各データベースはその一部のデータを格納した不完全なデータベースと考えることにより広域データベース検索を行う方法を提案した[宮崎96]。本稿では、この方法を実現するため関係代数を拡張し広域データベース検索処理の具体化を検討する。以下、2章で不完全データベースによる広域検索、3章で拡張関係、4章拡張関係代数、5章で広域検索への適用について述べる。

## 2 不完全データベースによる広域検索

広域検索ではどこにどんなデータがあるか正確にはわからない場合でも問合せを可能としたい。

このため関係データベースを拡張し不完全データベースの概念を導入する。不完全データベースでは分かっている情報のみを表現するので、タブルによって存在する属性が異なっても良い。また、問合せの条件にマッチするかどうか不明な時や条件付のデータを表現するために条件付タブルの概念を導入する。関係はこのように拡張されたタブルの集合である。不完全データベースによる広域検索では問合せは情報を持っている可能性のあるサイトに送られ問合せを受けたサイトは答えられるものだけ答える方式を取る。以下に例を挙げて説明する。

[例2. 1] 以下のスキーマのデータがあるとする。

学生 {名前、住所、電話}

また以下のようなタブルがある。

学生 {名前=千葉太郎、住所=習志野市、電話=0474-00-0000}

千葉太郎の年齢と住所に関する問合せ

? 学生 (名前=千葉太郎、年齢=X、住所=Y)を考える。従来のDBMSではこの問合せはスキーマと不一致のためエラーになる。不完全データベースでは存在する情報から、

学生 (名前=千葉太郎、住所=習志野市)を解とできる。年齢に関するデータはないのでこのような解を部分解と呼ぶ。これに対し問合せた全ての情報を含む解を完全解と呼ぶ。

[例2. 2] 複数のサイトに問合せを行う場合

? 学生 (名前=千葉太郎、住所=X、単位数=Y)を学生課と教務課に送りそれぞれ以下の部分解が得られたとする。

(学生課)

学生 (名前=千葉太郎、住所=習志野市)

(教務課)

学生 (名前=千葉太郎、単位数=80)

この 2 つの部分解の自然結合を取れば完全解が得られる。

### 3 拡張関係

不完全データベースに対する問合せでは、問合せた属性がデータベースに存在しないということが多々ある。また演算に関しては、広域で問合せを行う場合、各リレーションの属性がまちまちで、和両立が成立する可能性が低く複数サイトからの解の和がとれない場合がある。このような状況を解決するために関係を拡張する。

#### 3. 1 タプルの表現

レコード記法[YM 9 4]を用いてタプルの定義を行う。レコード記法では存在する属性のみ扱うことができる。

##### (1) 単純タプル

$$a(l_1=b_1, l_2=b_2, \dots, l_n=b_n)$$

ここでは  $a$  は関係名、  $l_i$  は属性名である。値の分からない属性は書かない。

##### (2) 条件付タプル

$$a(l_1=b_1, l_2=b_2, \dots, l_i=b_i) \text{ if } a(l_{i+1}=c_1, \dots, l_m=c_m)$$

これは  $\text{if}$  以下が真なら  $a(l_1=b_1, l_2=b_2, \dots, l_i=b_i)$  が真であることを表す。また、  $\text{if}$  以下が不等号条件を表すこともできる。

本稿では例に関してレコード記法の他に表によって表現する。表では存在する全ての属性を表記する。これ以降、レコード記法と従来の表形式を併用する。タプルによって存在しない属性は空欄とする。

#### 3. 2 拡張関係

拡張関係は同じ関係名の単純タプルと条件付タプルの集合として定義される。拡張関係に対して

以下の 2 つの概念を定義する。

##### (1) 共通属性制約：

・関係のすべてのタブルに少なくとも 1 つの共通属性が存在する。

共通属性制約は通常の関係の主キー制約に対応する。広域データベースで主キー制約を適用するのは困難である。したがって最小限の制約を導入した。

##### (2) 和両立

・2 つの関係に共通の属性があり、どのタブルにもその属性が存在する。

通常の関係の和両立の概念は制限が厳しいのでこのように変更した。

### 4 拡張関係代数

不完全データベースによる広域検索のために前節ではその基本的な部分の拡張を行ってきた。この節では、実現方法を検討するため、関係代数を拡張する。

#### 4. 1 集合演算

集合演算は従来の演算と同様に定義される。

和：二つの関係が和両立の時

$$A \cup B = \{x : x \in A \vee x \in B\}$$

と定義する。つまり  $A$  または  $B$  に属するタブルすべての集合。

差：二つの関係が和両立の時

$$A - B = \{x : x \in A \vee \neg(x \in B)\}$$

と定義する。つまり  $A$  に属しつつ  $B$  に属さないタブルの集合。

共通部分：二つの関係が和両立の時

$$A \cap B = \{x : x \in A \wedge x \in B\}$$

と定義する。つまり  $A$  と  $B$  に属すタブルの集合。

直積： $A$  と  $B$  を二つの拡張関係とする時

$$A \times B = \{(a, b) : a \in A \wedge b \in B\}$$

ここに  $a=(x_1, x_2, \dots, x_m), b=(y_1, y_2, \dots, y_n)$  とする時、 $(a, b)=(x_1, x_2, \dots, x_m, y_1, y_2, \dots, y_n)$  なる  $m+n$  項のタブルである。

[例 4. 1] 2つの拡張関係があるとする。

DB学会

AI学会

会員名	所属	連絡先
菅原	K55	2011
佐藤	K55	2035
森	K55	2038
鈴木	K81	2201

会員名	所属	住所
森	K55	町田
鈴木	K81	神戸
田中	K81	水戸

DB学会  $\cup$  AI学会

会員名	所属	連絡先	住所
菅原	K55	2011	
佐藤	K55	2035	
森	K55	2038	
鈴木	K81	2201	
森	K55		町田
鈴木	K81		神戸
田中	K81		水戸

## 4. 2 固有演算の拡張

ここでは固有演算の拡張の定義をする。

射影：

$X$ を属性の集合とする。関係Aの射影  $\pi_X(A)$  は、Aの各タブルから  $x$  に属する属性で属性値があるものだけを取り出したタブルの集合。また、 $X$ に属する属性が1つも存在しないタブルの射影は属性なしのタブルとする。

選択：

関係Aの選択  $\sigma_F(A)$  は、Aのタブルのうち条件Fを真とするまたは真とする可能性のある全てのタブルの集合である。対象とする属性値が存在しない場合、結果は条件付タブルとなり、対象とする属性が条件付の場合はタブルの条件の値が選択条件を満たせば結果は条件付タブルとなる。

$\theta$ -結合：

二つの関係AとBの  $\theta$ -結合  $A \theta B$  は  $\sigma_F(A \times B)$  で定義される。対象とする属性が存在しない場合条件付にして結合し、対象とする属性値が条件付の場合は条件の値を対象として結合する。

自然結合：

関係AとBの自然結合  $A \bowtie B$  は、AとBの直積から共通の属性の値が等しいものだけを選び、その結果から共通の属性の一方だけを残すように射影した関係である。一方の属性が存在せず条件部分に属性が現れる場合は条件を満たせば結合できる。

通常の関係代数では自然結合は  $\theta$ -結合と射影の組合せで表現できるが、拡張関係ではタブル対毎に共通属性が異なるので組合せでは表現できない。

無条件タブル選択：

関係Rのタブル内で条件付タブルを取り除く演算を無条件タブル選択と呼び  $\text{unc}(R)$  と書く。

## 4. 3 解のマージ

[例 2. 2] のように複数のサイトから部分解が得られた場合、それらをもとに完全解が生成できる場合がある。例えば複数のサイトから得られた解の和が以下のようになつたとする。

会員名	所属	連絡先	住所
菅原	K55	2011	
佐藤	K55	2035	
森	K55	2038	
鈴木	K81	2201	
田中	K81		水戸
森	K55		町田

このとき、ユーザが「会員名=森」が同じものだと考え、下記のようにまとめて一つのタブルにしたいとする。

会員名	所属	連絡先	住所
菅原	K55	2011	
佐藤	K55	2035	
森	K55	2038	町田
鈴木	K81	2201	
田中	K81		水戸

この演算を“マージ”と呼ぶ。マージは“縮退”演算と自然結合の組合せで定義できる。

#### [タブルの順序]

同じリレーションの2つのタブル $T_i$ と $T_j$ を考える。以下の2つの条件が成立つ時 $T_i < T_j$ と書く。  
“<”は半順序関係である。

(1)  $T_i$ の属性名がすべて $T_j$ に存在し、その属性値すべてが等しい。

(2)  $T_i$ に条件付属性が存在する時は、

(a) $T_j$ の属性がその条件を満たす。

または、

(b) $T_j$ の条件付属性の条件が $T_i$ の条件を含む。

#### [縮退]

関係 $R$ のタブルの中で $T_i < T_j$ となるタブル $T_j$ が存在するようなタブル $T_i$ をすべて消去する演算を縮退と呼び $red(R)$ と書く。縮退の結果は消去する順番に依存せず一意に定まる。

A1	B		
A1	B	C1	D
A2		C2	



A1	B	C1	D
A2		C2	

#### [マージ]

$$mrg(R) = red(R \bowtie R)$$

で定義される演算をマージと呼ぶ。すなわち同じ属性値のタブルの自然結合により情報を統合し、不要なタブルを消去するのがマージである。

## 5 広域検索への適用

問合せはレコード記法で以下のように表現する。

$$? a(l_1=b_1, l_2=b_2, \dots, l_i=b_i, l_{i+1}=X_1, \dots, l_m=X_m)$$

属性値の定数は条件を、変数は求める対象を表す。また、複数の拡張述語の論理積を書くことができる。

処理手順：

- ・問合せを各サイトへ送る。
- ・各サイトから解が帰ってくる。
- ・解の和をとる。
- ・場合によってユーザがマージ、無条件タブル選択を行う。

### 5. 1 適用例

#### [例 5. 1]

学生課

就職課

学生

学生

名前	住所
千葉太郎	千葉市
大塚真吾	習志野市

名前	電話番号
千葉太郎	00-0000
大塚真吾	11-1111

「大塚真吾の住所と電話番号」を問合せをする。

? 学生 (名前=大塚真吾、住所=Y、

電話番号=Z)

・拡張関係代数で表現すると

$\pi(\text{住所}, \text{電話番号})(\sigma \text{ 名前}=\text{大塚真吾}(\text{学生}))$ .

それぞれのサイトでこの問合せを処理し、その結果の和をとる。

名前	住所	電話番号
大塚真吾	習志野市	
大塚真吾		11-1111

これをマージすると

名前	住所	電話番号
大塚真吾	習志野市	11-1111

となる。

### [例 5. 2]

学生課

学生

名前	住所
千葉太郎	千葉市
大塚真吾	習志野市

教務課

学生

名前	単位数
千葉太郎	85
大塚真吾	123

「単位数が 100 以上の学生の名前と住所」を問合せする。

? 学生 (名前 = X、住所 = Y、

単位数 = Z) , Z > 100

・拡張関係代数では

$\pi$  {名前、住所、単位数} ( $\sigma$  単位数  $\geq 100$  (学生))

結果の和は

名前	住所	単位数
千葉太郎	千葉市	if 単位数 $\geq 100$
大塚真吾	習志野市	if 単位数 $\geq 100$
大塚真吾		123

となる。マージすると

名前	住所	単位数
千葉太郎	千葉市	if 単位数 $\geq 100$
大塚真吾	習志野市	123

となる。ここで無条件タブル選択を行うと、

名前	住所	単位数
大塚真吾	習志野市	123

となり解が得られる。

## 5. 2 問題点とその解決法

(1) [例 5. 2] で学生課へ問合せた解が全て条件付き解になってしまうという問題点がある。

これを解決するには以下のようない方法がある。

- ・問合せの選択の条件を 1 つも満たしていないものは可能性解としサイト名だけ返す。
- ・他のサイトに問合せた解に部分解がある場合、それを基に可能性解があるサイトに再度問合せを行う。

(2) 広域検索を行うときユーザは問合せをするサイトがどのような情報を持っているかを把握するのが困難である。問合せに SQL の \* (全属性指定) に相当する表現を導入すればユーザは問合せ先にある情報をすべて引き出す事ができ、その情報を有効に活用できる。

## 6 まとめ

広域データベースでは従来のマルチデータベースの方法のようにあらかじめスキーマ統合を行っておくことが困難である。スキーマ統合を行った場合でも、情報のミスマッチをいかに解消するかの問題がある。本稿では広域ネットワークに散在するデータベースに対し、問合せを行う方法を具体化するため関係代数の拡張を提案した。実現には多くの研究課題があるが広域データベース検索が可能になればネットワークの情報利用の可能性は飛躍的に高まるであろう。今後はリレーションナルデータベースシステムを利用したプロトタイプを作成し、その有効性を検討していく予定である。

## 参考文献

- [宮崎 96] 宮崎収兄, “不完全データベースと広域データベース検索” 情報処理学会データベースシステム研究報告, 96-DBS-106, pp. 131-138  
[YM 94] 横田一正, 宮崎収兄, “新データベース論” 共立出版, 1994