

関係データベースにおける

Representative Instance に関する一結果

阪大 基礎工 岩崎元昭

伊藤 実

高 忠雄

1. まえがき

最近、universal instance を一般化した概念として representative instance が提案された⁽⁴⁾⁽⁶⁾⁽⁸⁾。representative instance は、データベースに属する各関係が universal instance の射影でない場合においてもデータベース全体の内容を表現できることひとつとして提案され、この関係がデータベーススキームで指定されたすべての制約を満たさなければならぬと仮定されている。本論文では、制約として関数従属だけを考え、次の 2 つの多項式時間手続きを示した。

(1) 与えられたデータベーススキームが整合性を保つかどうか、すなわち、与えられたデータベーススキームの任意のデータベースに対して、その representative instance が常にすべての関数従属を満たすかどうかを判定する。

(2) データベーススキームが整合性を保つとき、represen-

tative instance の全射影⁽⁶⁾を効率よく計算するための関係表現を構成する。

これらの結果は、文献(7)の結果をさらに一般化したものである。

2. 諸定義

$R = \{A_1, \dots, A_n\}$ を属性の集合とし、各属性 A_i の定義域を $\text{dom}(A_i)$ とする。直積 $\text{dom}(A_1) \times \dots \times \text{dom}(A_n)$ の有限な部分集合 r を R 上の関係という。 t を $\text{dom}(A_1) \times \dots \times \text{dom}(A_n)$ の元（組）といふ）とし、 A を R に属する属性とすると、 t の A 成分を $t[A]$ と表す。同様に、 R の部分集合 $X = \{A'_1, \dots, A'_m\}$ に対して、 $(t[A'_1], \dots, t[A'_m])$ を $t[X]$ と表す。以下では、属性を表すのに A, B, C, \dots 等の文字を用いて、属性の集合を表すのに Z, Y, X, \dots 等の文字を用いる。属性集合 X, Y に対して、その和集合 $X \cup Y$ を XY と略記する二ことがある。また、ひとつ属性 E からなる集合 $\{A\}$ を単に A と書くことがある。

R 上で定義される関数従属(FD) とは、 $X \rightarrow A$ なる文である。但し、 $X \subseteq R$, $A \in R$ である。 R 上の関係 Y に属する任意の二つの組 t, s に対して、もし $t[X] = s[X]$ ならば $t[A] = s[A]$ であるとき、 Y は $X \rightarrow A$ を満足すといふ。FD の集合

$\{X \rightarrow A_1, \dots, X \rightarrow A_i\}$ をまとめて $X \rightarrow A_1 \dots A_i$ と書くことがある。FDのある集合Fに対して、Fに属するすべてのFDを満たす任意の関係において常に成立する(F以外の)FDが存在する。そのようなすべてのFDの集合を(Fも含めて) F^+ と表す。Rの部分集合Xに対して、Fに属する閉包(closure)を $\text{closure}(X, F) = \{A \mid X \rightarrow A \in F^+\}$ と定義する。

$\text{closure}(X, F)$ を求める $O(|F|)$ 時間の手続きが知られてる⁽¹⁾。

本論文では、一般に変数を含む関係を考える。すなはち、関係に属するある組のある成分が定数ではなく変数であるのも許す。2つの組 t_1, t_2 に対して、 $t_1[A] = t_2[A]$ であるのは、 t_1 と t_2 が A 成分に同じ定数または同じ変数をもつ場合である。 γ を一般に変数を含む R 上の関係とし、 $F \in R$ 上で定義される FD の集合とする。 γ に対する F のまととの chase⁽⁵⁾ とは、Fに属する各 FD に対する次に定義する FD 規則を適用して γ に現れる変数を定数または変数に置き換える操作である。

[FD規則] $X \rightarrow A \in F$ に属する FD とし、 t_1, t_2 を $t_1[X] = t_2[X]$ かつ $t_1[A] \neq t_2[A]$ なる γ に属する組とする。 t_1 は A 成分に変数ひをもつとする。 $X \rightarrow A$ に関する FD 規則によって、 γ の A 成分に現れる変数ひはすべて $t_2[A]$ に置き換えられる。このとき、 $t_1[A], t_2[A]$ とも異なる定数をも

つ場合には、FD規則は適用されず、Yに対するFのまととのchaseはFを満たさないという。

ある属性集合U上のデータベーススキームとは、2字組の集合 $R = \{ \langle R_1, F_1 \rangle, \dots, \langle R_n, F_n \rangle \}$ である。但し、 $U = R_1 \cup \dots \cup R_n$ かつ各 F_i は R_i 上で定義されるFDの集合である。各 R_i を関係スキームという。関係の集合 $I = \{ Y_1, \dots, Y_n \}$ に対して、もし各 Y_i が R_i 上の関係であり、かつ、 F_i を満たすならば、IをRのデータベースという。すなわち、IはデータベーススキームIRのある時点での値を表す。Iに属する関係には変数は現れないものとする。以下では、データベーススキーム $R = \{ \langle R_1, F_1 \rangle, \dots, \langle R_n, F_n \rangle \}$ が与えられたとき、 $U = R_1 \cup \dots \cup R_n$ 、 $F = F_1 \cup \dots \cup F_n$ と表し、次の2つの仮定をみく。

1. (universal relation scheme assumption) $\forall (X \rightarrow A \in F^+ \text{ かつ } X \subseteq R_i) \Rightarrow X \rightarrow A \in F_i^+$ である。この仮定は、各属性は下下一つの役割しか果たさない。すなわち、任意の属性集合Xに対して、Xの属性間の関係(relationship)は高々一つしか存在しないことを表す。これは関係データベースの設計理論に関する多くの論文で仮定されており。

2. (FDの集合の最適化) 各 F_i に属する任意のFD $X \rightarrow A$ に対して、

$$(a) \quad X \rightarrow A \notin (F_i - \{X \rightarrow A\})^+$$

(b) 任意の $X' \subseteq X$ に対して、 $X' \rightarrow A \notin F_i^+$

が成立する。与えられた F_i が (a), (b) を満足しない場合、 $O(|F_i| |F_j|)$ 時間でこれらの性質を満足すように変換できる⁽²⁾。

本論文では、関係に対する演算として射影、結合、及び、関係和を考える。R 上の関係 Y 及び R の部分集合 X に対して、Y の X 上への射影を、 $Y[X] = \{t[X] \mid t \in Y\}$ と定義する。R_i 上の関係 Y_i 及び R_j 上の関係 Y_j に対して、Y_i と Y_j の結合を、 $Y_i \bowtie Y_j = \{t \mid t[R_i] \in Y_i \text{ かつ } t[R_j] \in Y_j\}$ と定義する。R 上の 2 つの関係 Y_i, Y_j に対して、Y_i と Y_j の関係和を、 $Y_i \cup Y_j = \{t \mid t \in Y_i \text{ または } t \in Y_j\}$ と定義する。R 上の関係表現とは、R に属する関係スキーム R₁, …, R_n を変数とし、射影、結合、関係和を基本演算とする式で、次のように定義される。

(1) R_i はそれ自身で関係表現である。

(2) E₁, E₂ を関係表現とすると、E₁[X], E₁ \bowtie E₂, E₁ \cup E₂ は関係表である。

I = {Y₁, …, Y_n} を R のデータベースとすると、E は I から一つの関係 E(I) への写像を定義する。E(I) の値は、E 中に現れた各変数 R_i に関係 Y_i を値として代入し、基本演算の定義に従って計算することにより求まる。R_i 上の関係 Y_i 及び R_j 上の関係 Y_j に対して、Y_j が R_i \cap R_j \rightarrow R_i なる FD を満足すなら

ば、 $R_1 \bowtie R_2$ を拡張結合 (extension join) という⁽³⁾。

$I = \{R_1, \dots, R_n\}$ を R のデータベースとする。 R_i に属する各組の $U - R_i$ 成分に相異なる変数を追加することにより、 R_i は U 上の関係とみなせる。 R_i の U 上への関係への拡張を、

$\text{aug}_U(R_i) = \{t \mid t[R_i] \in R_i \text{かつ } U - R_i \text{ に属する各属性 } A \text{ に対して、 } t \text{ の } A \text{ 成分は他の } t' \text{ に現れない変数}\}$

と定義する。更に、 $\text{aug}_U(I) = \bigcup_{i=1}^n \text{aug}_U(R_i)$ と定義する。このとき、 I の表現例 (representative instance) とは、 $\text{aug}(I)$ に対する F のもとでの chase によって最終的に得られる関係のことをあり、 $\text{rep}(I)$ と書く。

R の任意のデータベース $I = \{R_1, \dots, R_n\}$ に対して、その表現例 $\text{rep}(I)$ が常に F を満たすならば、 R は整合性を保つといふ。従って、 R が整合性を満たせば、 R の任意のデータベースに対して、 $\text{aug}_U(I)$ に対する F のもとでの chase は必ず F を満たす。

3. 整合性判定手続き

本節では、データベーススキーム R が任意に与えられたときに、 R が整合性を保つかどうかを判定する手続きを示す。

また、手続きの時間複雑度の評価を行う。但し、手続きの正当性の証明は省略する。

[整合性判定手続き]

入力：データベーススキーム $R = \{<R_1, F_1>, \dots, <R_n, F_n>\}$

方法：ある関係スキーム R_i に対して次の手続き EXAM(R_i) が
“no”を返せば R は整合性を保たない。すべての R_i に対して
EXAM(R_i) が “yes” を返せば R は整合性を保つ。

procedure EXAM(R_i)

begin

a. i 以外の各 j ($1 \leq j \leq n$) に対して、 $G_j \leftarrow F_j$, $H_j \leftarrow \emptyset$ と
する。 $(G_j$ は F_j に属する FD でまだ調べていなきものの、 H_j は
既に調べられたものの集合をそれぞれ表す。)

b. $S \leftarrow R_i$ (以下では、 S に対して $R_i \rightarrow S \in F^+$ が成立する。)

c. ある G_j ($1 \leq j \leq n$, $j \neq i$) に $X \subseteq S$ なる FD $X \rightarrow A$ が存
在する限り、次の 1~2 を実行する。

1. (G_j に $X \subseteq S$ なる FD $X \rightarrow A$ が存在するとて)

G_j から

(i) $X_1 \subseteq S, \dots, X_p \subseteq S$

(ii) $\text{closure}(X_1, F_j) = \dots = \text{closure}(X_p, F_j)$

(iii) $\text{closure}(X_1, F_j)$ は G_j に属する FD の中で極小 ($X \subseteq S$ な
る任意の FD $X \rightarrow A \in G_j$ に対して $\text{closure}(X, F_j) \neq \text{closure}(X_1, F_j)$)
を満たす FD の集合 $X_1 \rightarrow A_1, \dots, X_p \rightarrow A_p$ を任意に選ぶ。

2. 上記 1. で選ばれた各 FD $X_g \rightarrow A_g$ に対して次の条件 (不

整合条件といふ)を満たすかどうかを調べる。

[不整合条件] $S \cap \text{closure}(X_g, F_j) - \text{closure}(X_g, H_j) \neq \emptyset$
 もしくは、不整合条件を満たすものが存在すれば“no”を返す。
 もなければ、

$$S \leftarrow S \cup \text{closure}(X_i, F_j)$$

$$H_j \leftarrow H_j \cup \{ X \rightarrow A \mid X \rightarrow A \in G_j \text{ かつ } X \subseteq \text{closure}(X_i, F_j) \}$$

$$G_j \leftarrow G_j - H_j$$

とする。

d. ($C \geq "no"$ が返されずに終了した場合) “yes”を返す。

end EXAM \square

[定理1] データベーススキーム $R = \{ \langle R_1, F_1 \rangle, \dots, \langle R_n, F_n \rangle \}$
 が整合性を保つかどうかは、 $O(n|F| |F|)$ 時間で判定できる。
 但し、 $F = F_1 \cup \dots \cup F_n$ とする。 \square

4. 表現例の全射影の計算

$R = \{ \langle R_1, F_1 \rangle, \dots, \langle R_n, F_n \rangle \}$ を整合性を保つデータベーススキームとし、 $I = \{ r_1, \dots, r_n \}$ を R のデータベースとする。
 利用者によるデータベースに対する質問が属性集合 I を参照するとすると、その答を求めるためには表現例の上への射影を求めることが必要である。利用者は成分に変数を含まない組だけに興味があると仮定する。一般に変数を含む関係

γ の Z 上への全射影を $\gamma[Z\text{-total}] = \{t[Z] \mid t \in \gamma \text{ かつ } t \text{ は } Z \text{ 成分に変数を含まない}\}$ と定義する。 $\text{rep}(I)[Z\text{-total}]$ を求めるのに、 $\text{rep}(I)$ を求めた後 Z 上への射影をとるには多くの時間を必要とする効率的ではない。本節では、 $R_1 \cup \dots \cup R_n$ の任意の部分集合 Z が与えられたとき、 R の任意のデータベース I に対して $E_Z(I) = \text{rep}(I)[Z\text{-total}]$ を満たす関係表現 E を多項式時間で求める方法を示す。 E_Z は $E_1 \cup \dots \cup E_n$ の形をしており、各 E_i は拡張結合の系列になっている。従って、 $E_Z(I)$ は効率よく求めることができることができる。

$E_Z = E_1 \cup \dots \cup E_n$ の求め方を簡単に説明すると、各 E_i は EXAM(R_i) において、 $Z \subseteq S$ となるような最短の FD の並び方を。 $X_1 \rightarrow A_1 \in F_{j_1}, \dots, X_m \rightarrow A_m \in F_{j_m}$ とする。このとき、 $E_i = R_i \cup R_{j_1}[\text{closure}(X_1, F_{j_1})] \cup \dots \cup R_{j_m}[\text{closure}(X_m, F_{j_m})]$ とする。但し、 $\text{closure}(X_1, F_{j_1}) \subseteq R_{j_1}$ である。もし、 $Z \neq \text{closure}(R_i, F_1 \cup \dots \cup F_n)$ ならば E_i は定義されない。

[定理 2] R を整合性を保つデータベーススキームとする。属性集合 Z が任意に与えられたとき、 R の任意のデータベース I に対して $\text{rep}(I)[Z\text{-total}] = E_Z(I)$ となる関係表現 E_Z は $O(n|F| \|F\|)$ 時間で求まる。□

参考文献

- (1) Beeri, C. and P. A. Bernstein, "Computational problems related

- to the design of normal form relational schemes," ACM Trans. on Database Systems, Vol.4, No.1 (March 1979), pp. 30-59.
- (2) Bernstein, P. A., "Synthesizing third normal form relations from functional dependencies," ACM Tran. on Database Systems, Vol. 1, No. 4 (Dec. 1976), pp. 277-298.
- (3) Honeyman, P., "Extension joins," Proc. Int. Conf. on VLDB, 1980, pp. 239-244.
- (4) Honeyman, P., "Testing satisfaction of functional dependencies," Proc. XPI Conf., Stonybrook, N.Y., June, 1980.
- (5) Maier, D., A.O. Mendelzon and Y. Sagiv, "Testing Implications of dependencies," ACM Tran. on Database Systems, Vol.4, No. 4 (Dec. 1979), pp. 455-469.
- (6) Sagiv, Y., "Can we use the universal instance assumption without using nulls?" Proc. ACM-SIGMOD Int. Conf. on Management of Data, Ann Arbor, April 1981, pp. 108-120.
- (7) Sagiv, Y., "A characterization of globally consistent databases and their correct access paths."
- (8) Vassiliou, Y., "A formal treatment of imperfect information in database management," Technical Report CSRG-123, University of Toronto, Nov., 1980.