

# 統計データベースのある最適設計手法について

小林康幸

島根大学

素データから作成される要約表の集まりが与えられた時、集約、射影による操作によってそれらの要約表をすべて導出できる要約表の集まりの中で、総レコード数が最小となるものを見つける問題を考える。与えられた要約表は、利用者が必要とする要約表であり、求める要約表は、計算機にデータベースとして格納する要約表（ベース表）である。すなわちこの問題は、ある意味での計算機での格納容量を最小にする、データベース設計の問題である。

ここでは、まずこの問題を数学的に定式化し、問題を解くためのいくつかの定理を述べた後、これらの定理を用いて作った、問題の解を求めるための方法を与える。

## AN OPTIMAL DESIGN METHOD OF STATISTICAL DATABASES

Yasuyuki Kobayashi

Shimane University

1060, Nishikawatsu-cho, Matsue, 690, Japan

This paper treats a problem on statistical database designs; for a given set of summary tables, find a set of base tables to be organized as a database such that all given summary tables are derivable from them and that the total number of records is minimum among such sets of base tables.

After formalizing this problem and describing some theorems to solve the problem, we give a method to obtain a solution of the problem.

## 1. まえがき

統計データは、通常大容量であることおよび、個々のレコードの機密保護の点から、素データを計算機に格納せず、素データを集約したいいくつかの要約表を格納することが多い。この場合、利用者が必要とする要約表の集まりが与えられた時、① どのような要約表の集まりを計算機に格納すれば、それらの要約表を選択、射影、結合、集約などの操作によって導出できるか、② ①で得られた要約表の集まりの中で、計算機での格納容量が最小となるものはどれか、という問題が生じる。

ここでは、これらの問題の定式化および、問題を解くための方法について述べる。

## 2. 問題の定式化

工業統計や商業統計のように、通常統計表は、一つの調査によって得られた素データ（一次統計）から、集約によつていくつかの要約表（二次統計）が作成される。利用者は、与えられた要約表の集まりから、各自に必要な表を、選択、射影、結合、集約などの操作で導出することになる。

素データ  $R$  を領域の集まり  $\{D_i : i = 1, \dots, N\}$  の直積  $D = \prod_{i=1}^N D_i$  の有限部分集合とする。また

任意の集合  $A$  に対し、 $A^*$  を  $A$  の有限個の元の列のすべての集まりとし、 $A$  上の要約演算子  $g$  を、 $A^*$  の元  $a$  の任意の置換  $a'$  に対し  $g(a') = g(a)$  となる関数と定義する。さらに  $A$  上の要約演算子  $g$  が次の条件をみたす時、 $g$  を結合的要約演算子と呼ぶ：

$A^*$  の任意の元  $a, a'$  に対し、 $g(g(a), g(a')) = g(a \cup a')$ ,

但し  $a \cup a'$  は  $a$  および  $a'$  を一列に並べた  $A$  の列とする。合計、最大値などは結合的要約演算子であるが、中央値は結合的でない要約演算子である。

素データ  $R$  ( $\subset \prod_{i=1}^N D_i$ )、 $\{1, \dots, N\}$  の空でない真部分集合  $X$  および  $\{1, \dots, N\} - X$  の元

$j$  に対し、 $D_X$  から  $D_j$  への関数  $E_{j,x}$  を

$E_{j,x}(d) = (p_{r_j}(r) : r \in R, p_{r_j}(r) = d) \quad (d \in D_X)$

と定義する（但し  $\{1, \dots, N\}$  の空でない任意の部分集合  $X$  に対し、 $D_X = \prod_{i \in X} D_i$  および  $D_X$

から  $D_i$  への射影を  $p_{r_i}$  とする。）。

また、 $\emptyset \neq Y = Y(G) (\subset \{1, \dots, N\} - X)$  に対し、 $G = \{g_j : j \in Y(G)\}$  を  $D_j$  上の要約演算子  $g_j$  の集合とする時、 $D_X$  から  $D_{X \cup Y}$  への関数  $S_{G,X}$  、

$S_{G,X}(d) = (d, (g_j(E_{j,x}(d)) : j \in Y)) \quad (d \in D_X)$

を  $R$  上の集約関数と呼び、 $X \subset X' \subset \{1, \dots, N\}$  なる  $X'$  に対し、 $D_{X'}$  から  $D_{X \cup Y}$  への関数  $S_{G,X} \circ p_{r_X}$  を同じ記号  $S_{G,X}$  で表すことにする（ $X' \supset X$  の時、 $D_{X'}$  から  $D_X$  への射影  $p_{r_X}$  は、 $X$  の任意の元  $i$  に対し、 $p_{r_X} \circ p_{r_X}$  が  $D_{X'}$  から  $D_i$  への射影  $p_{r_i}$  と一致する関数とする。）。 $S_{G,X}$  による  $R$  の像  $S(G, X; R) = S_{G,X}(R) (\subset D_{X \cup Y})$  を  $S_{G,X}$  による  $R$  の要約表と呼び、 $X$  および  $Y$  はそれぞれ分類項目および集約項目の集合と呼ぶ。

例 2.1.  $D_1 = \{ABC, ACC, NSS, QQQ, WSE, BTT\}$  ,

$D_2 = \{KYOTO, TOKYO\}$  とし、 $D_3, D_4, D_5, D_6$  を非負整数の集合とすると、表 2.1 の  $R$

( $\subset \prod_{i=1}^6 D_i$ ) は素データである。また  $X = \{2\}$  ,  $i = 5$  とすると、

$E_{i,x,R}(KYOTO) = (200, 310, 140, 100)$ ,  
 $E_{i,x,R}(TOKYO) = (580, 120, 230, 130)$ ,  
 $E_{i,x,R}(D_2) = \{(200, 310, 140, 100), (580, 120, 230, 130)\}$   
 となる。さらに  $X = \{2\}$ ,  $Y = \{5, 6\}$ ,  $g_5$  および  $g_6$  を SUM (合計) とすると、  
 $S_{G,x}(R) = \{(KYOTO, 750, 134), (TOKYO, 1060, 148)\}$   
 となる。

		表2.1 素データ R			
A	B	ABC	KYOTO	1	1 200 27
A	C	ACC	KYOTO	1	2 310 31
N	S	NSS	KYOTO	2	1 140 56
P	C	PCC	KYOTO	2	2 100 20
Q	Q	QQQ	TOKYO	2	2 580 21
W	S	WSE	TOKYO	2	1 120 74
B	T	BTT	TOKYO	1	2 230 35
A	X	AXT	TOKYO	1	1 130 18

素データ R ( $\subset \prod_{i=1}^N D_i$ ) および R 上の集約関数の有限集合の集まり F S A F が与えられた時、F S A F の元  $S$  および  $S'$  が次の条件 (C 1) を満たす時、 $S$  は  $S'$  から導出可能 ( $S' \rightarrow S$  と書く) であると呼ぶ：

(C 1)  $S$  の任意の元  $S_{G,x}$  に対し、

- i)  $\cup \{G': S_{G',x'} \in S'\} = G$  かつ
- ii)  $S'$  の任意の元  $S_{G',x'}$  に対し

$$X' \supset X \text{ かつ } S(G', X; S(G', X'; R)) = S(G', X; R)$$

を満たす  $S'$  の部分集合  $S'$  が存在する (但し上の等式で、 $S(G', X; R')$  は、 $R' = S(G', X'; R)$  ( $\subset \prod_{j \in X' \cup Y} D_j$ ) を素データとした時の  $S_{G',x'}$  による  $R'$  の要約表である。)。

例 2.2. 例 2.1において  $X_1 = \{2, 3\}$ ,  $X_2 = \{2, 4\}$ ,  $G_1 = \{g_5\}$ ,  $G_2 = \{g_6\}$ ,  $S' = \{S_{G_1,x_1}, S_{G_2,x_2}\}$ ,  $X = \{2\}$ ,  $G = \{g_5, g_6\}$ ,  $S = \{S_{G,x}\}$  とすると  $S_{G_1,x_1}(R)$ ,  $S_{G_2,x_2}(R)$ ,  $S_{G,x}(R)$  は次のようになり、 $S' \rightarrow S$  を容易に示すことができる。

$S_{G_1,x_1}(R)$	$S_{G_2,x_2}(R)$	$S_{G,x}(R)$
KYOTO 1 510	KYOTO 1 83	KYOTO 750 134
KYOTO 2 240	KYOTO 2 51	TOKYO 1060 148
TOKYO 1 360	TOKYO 1 92	
TOKYO 2 700	TOKYO 2 56	

F S A F の元  $S$  に対し、次の式で  $S$  のレコード数 N R e c ( $S$ ) を定義する：

$$NRec(S) = \sum_{S_{G,x} \in S} |S_{G,x}(R)|,$$

但し、集合Aに対し、 $|A|$ をAの元の数とする。

以上の定義を用いて、当初の目的であった”与えられた要約表をすべて導出できる要約表の集合の中で、計算機での格納容量の最小となるもの”的定義をする。

FSAFの元 $S_0$ が与えられた時、 $S \in FSAF$  が次の条件を満たす時、 $S$ は $S_0$ に対するレコード最小集合であると呼ぶ：

- i)  $S \rightarrow S_0$
- ii)  $NRec(S) = \min \{NRec(S') : S' \in FSAF \text{かつ } S' \rightarrow S_0\}$

またこの時、 $S$ は $S_0$ に対する統計データベース設計の問題に対する解あるいは単に $SDD(S_0)$ の解と呼ぶ。

### 3. 問題の解法

ここでは、まず次のような問題を考える。

空でない有限集合U、Uの被覆B(BはUの部分集合の集まりで、 $\langle B \rangle$  (集合族Aに対し、 $\langle A \rangle = \bigcup_{A \in A} A$ とする。) およびUの全ての元xに対する非負整数f(x)が与えられた時、

次の条件を満たす $2^U - \{\emptyset\}$  の部分集合Aを $\langle U, B, f \rangle$ に対する最小積和分解問題の解あるいは $\langle U, B, f \rangle$ に対するMSPDと呼ぶ：

- i)  $A \succ B$ 、すなわちBの任意の元Bに対し、Aの元Aが存在し、 $A \supseteq B$
- ii)  $A' \succ B$ となる任意の $2^U - \{\emptyset\}$  の部分集合A'に対し、 $SP_f(A) \leq SP_f(A')$   
但し $2^U - \{\emptyset\}$  の任意の部分集合A'に対し、 $SP_f(A') = \sum_{A \in A'} \prod_{x \in A} f(x)$ と定義する。

すると次のようなことが示せる。

定理1.  $R(\subset \prod_{i=1}^N D_i)$ を素データ $S_0(\in FSAF)$ をR上の有限個の集約関数の集まりとし、

$B = A(S_0) = \{X : S_{G,x} \in S_0\}$  および $U = \langle B \rangle$ とおく。 $pr_U(R) = \prod_{i \in U} pr_i(R)$ および

要約演算子の有限集合 $G_0$ が存在し、 $S_0$ の任意の元 $S_{G,x}$ に対し、 $G = G_0$ であることを仮定する。

- i)  $S$ が $SDD(S_0)$ の解なら、 $A(S) \cap U = \{X \cap U : S_{G,x} \in S\} - \{\emptyset\}$  は $\langle U, B, f \rangle$ のMSPDである(但し、 $f(i) = |pr_i(R)|$  ( $i \in U$ ) )。
- ii) 逆に、Aが $\langle U, B, f \rangle$ の解なら、任意の要約演算子の有限集合 $G_0$ に対し、 $S(A) = \{S_{G,x} : X \in A\}$  は $SDD(S_0)$ の解である。

例3.1. 例2.1において $Y = \{5, 6\}$ ,  $U = \{2, 3, 4\}$ ,  $f_5 = f_6 = SUM$ ,  $G = \{g_5, g_6\}$ ,  $f(i) = |pr_i(R)| = 2$  ( $2 \leq i \leq 4$ ),  $B = \{(2, 3), (2, 4), (3, 4)\}$ ,  $S_0 = \{S_{G,B} : B \in B\}$  とすると、 $pr_U(R) = \prod_{i \in U} pr_i(R)$ であることは容易にわかる。従って

$SDD(S_0)$ の解を求めるることは、次のように表わされた $\langle U, B, f \rangle$ に対するMSPDを求めることがある。

2 3 4

B <sub>1</sub>	1	1
B <sub>2</sub>	1	1
B <sub>3</sub>	1	1

f 2 2 2

さらに  $\{U\}$  が  $\langle U, B, f \rangle$  に対する M SPD であること、従って  $\{S_{G,U}\}$  が SDD( $S_0$ ) の解であることも容易にわかる。

従って、以降では統計データベース設計の問題を解くために、 $\langle U, B, f \rangle$  に対する M SPD を求め  
る方法を考える。

まず簡単な場合について考える。命題および定理の証明は、文献 [2] を参照されたい。

命題1.  $f(x) = 0$  となる  $U$  の元  $x$  が存在するなら、任意の M SPD  $A$  に対し、 $SP_f(A) = 0$  となる。

命題2.  $B \subseteq B'$  となる  $B$  の元  $B$  および  $B'$  が存在するなら、 $A$  が  $\langle U, B, f \rangle$  に対する M SPD であるための必要十分条件は、 $A$  が  $\langle U, B - \{B\}, f \rangle$  に対する M SPD であることである。

命題3.  $f^{-1}(0) = \emptyset$ ,  $f^{-1}(1) \neq U$  の時

$$V = f^{-1}(1), U' = U - V, f' = f|_U$$

とし、 $2^U - \{\emptyset\}$  の部分集合  $A$  に対し、 $A - V = \{A - V : A \in A\} - \{\emptyset\}$  と定義すると、

- i)  $A$  が  $\langle U, B, f \rangle$  に対する M SPD なら、 $A - V$  は  $\langle U', B - V, f' \rangle$  に対する M SPD となる。
- ii) 逆に  $A'$  が  $\langle U', B - V, f' \rangle$  に対する M SPD なら、 $\{A' \cup V : A' \in A\}$  は  $\langle U, B, f \rangle$  に対する M SPD となる。

例3.2. 図3.1に示される  $\langle U, B, f \rangle$  に対する M SPD を求めるためには、命題3より図3.2に示される  $\langle U', B', f' \rangle$  に対する M SPD を求めれば良い。さらに命題2より、図3.3に示される  $\langle U', B' - \{B_4\}, f' \rangle$  に対する M SPD を求めれば良いことがわかり、 $\{B_1, B_2, B_3\}$  が M SPD であることが簡単に求められる。

	u <sub>1</sub>	u <sub>2</sub>	u <sub>3</sub>	u <sub>4</sub>
B <sub>1</sub>	1	1		1
B <sub>2</sub>	1		1	
B <sub>3</sub>		1	1	
B <sub>4</sub>	1			1
B <sub>5</sub>		1	1	1

f 3 2 2 1

図 3.1.  $\langle U, B, f \rangle$

	u <sub>1</sub>	u <sub>2</sub>	u <sub>3</sub>
B <sub>1</sub>	1	1	
B <sub>2</sub>	1		1
B <sub>3</sub>		1	1
B <sub>4</sub>	1		

f 3 2 2

図 3.2.  $\langle U', B', f' \rangle$

	$u_1$	$u_2$	$u_3$
$B_1$	1	1	
$B_2$	1		1
$B_3$		1	1
$f$	3	2	2

図 3.3.  $\langle U, B - \{B_4\}, f' \rangle$

命題1～3より、今後は次のことを仮定する：

(A1)  $B$ の任意の元 $B$ は $B - \{B\}$ の任意の元 $B'$ を含まない。

(A2)  $U$ の任意の元 $x$ に対し $f(x) \geq 2$ 。

上の仮定の下で、 $\langle U, B, f \rangle$ に対する特殊なタイプのMS PDを求めるために、いくつかの定義をする。 $2^U - \{\emptyset\}$ の部分集合 $C$ が次の条件を満たす時、 $C$ を分解不能族と呼ぶ：

(C2)  $\langle C \rangle$ の任意の元 $a, b$ は $C$ のある元 $c$ に含まれる。

仮定(A1)を満たす $2^U - \{\emptyset\}$ の部分集合 $B$ 、 $U$ の部分集合 $C$ および $\langle A \rangle \cap C = \emptyset$ となる $B$ の部分集合 $A$ に対し、 $C$ が次の条件を満たす時、 $C$ を $B$ における $A$ の $C$ -極大分解不能族あるいはC-MIFと呼ぶ：

(C3)  $A \subset C \subset B$ ,  $\langle C \rangle \cap C = \emptyset$ , かつ $C$ は分解不能族である。

(C4)  $A \subset C' \subset B$ ,  $\langle C' \rangle \cap C = \emptyset$ となる任意の $C'$ に対し、 $C' \supset C$ かつ $C'$ が分解不能族ならば、 $C' = C$ となる。

$A \subset B$ の時、 $C$ が次の条件を満たす時、 $C$ を $A$ における最大分解不能族と呼ぶ：

(C5)  $C \subset A$ かつ $C$ は分解不能族である。

(C6)  $C' \subset A$ かつ $\langle C' \rangle = \langle C \rangle$ ならば、 $C' \subset C$ である。

例3.3において、 $\{B_1, B_2, B_3\}$ および $\{B_1, B_4, B_5\}$ は、 $B$ における $\{B_1\}$ のO-MIFであり、従つて分解不能族でもある。また $u_3, u_4$ を含む $B$ の元が存在しないことから $B$ は分解不能族ではない。

### 例3.3.

$u_1 \quad u_2 \quad u_3 \quad u_4$

$B_1$	1	1	
$B_2$	1		1
$B_3$		1	1
$B_4$	1		1
$B_5$		1	1

定理2. (A1～2)の仮定の下で、条件(C7～9)を満たす $B$ の分割 $\beta$ に対する $A_\beta = \{\langle C \rangle : C \in \beta\}$ の集合の中に、 $\langle U, B, f \rangle$ に対するMS PDが存在する。

(C7)  $\beta$ の任意の元 $C$ は、分解不能族である。

(C8)  $\beta$ の任意の異なる元 $C, C'$ に対し、 $A \subset C$ かつ $\langle A \rangle \subset \langle C' \rangle$ なら、 $\langle C - A \rangle = \langle C \rangle$ かつ $C - A$ は分解不能族である。

(C9)  $|C| \geq 3$ となる $\beta$ の元 $C$ に対し、 $SP_r(\{C\}) < SP_r(C)$ が満たされる。

$C - M I F$ の定義における $B, A$ および $C$ に対し、

$$D^0(A, C) = D^0(C) = \{A \in B : A \cap C = \emptyset\}$$

$$D^{i+1}(A, C) = \{A \in D^i(A, C) : \langle A \rangle \text{の任意の元 } u \text{ に対し}, A \subset \langle D^i(A, C)_u \rangle\}$$

とすると、ある整数 $i$ が存在し、 $D^i(A, C) = D^{i+1}(A, C)$ となる。この時の $D^i(A, C)$ を $D^*(A, C)$ と表わす。次の定理により $D^*(A, C)$ を用いて $A$ の $C - M I F$ を求めることができる（但し、 $2^U - \{\emptyset\}$ の部分集合 $C, U$ の元 $u$ および $U$ の部分集合 $B$ に対し、 $C_u = \{A \in C : A \ni u\}$ と定義する。）。

定理3.  $A \not\subset D^*(A, C)$ 、または $D^*(A, C) = A$ かつ $A$ が分解可能族なら、 $A$ の $C - M I F$ は存在しない。

定理4.  $D^*(A, C) \supset A$ かつ $D^*(A, C)$ が分解不能族なら、 $D^*(A, C)$ は唯一の $A$ の $C - M I F$ である。

定理5.  $D^*(A, C) \supset A$ かつ $D^*(A, C)$ が分解可能族の時は、

1)  $B \cap C = \emptyset$ となる $B - A$ の元 $B$ が存在して、 $C$ が $A \cup \{B\}$ の $C - M I F$ ならば、 $C$ は $C \neq A$ となる $A$ の $C - M I F$ である。

2) 1) の仮定を満たす $C$ が存在しないなら、 $A$ の $C - M I F$ は、 $A$ が分解不能族なら $A$ となり、 $A$ が分解可能族なら存在しない。

定理6. 定理3～5を用いることにより、 $A$ のすべての $C - M I F$ を求めることができる。

次に $A$ の最大分解不能族を求めるための定理を述べる。

定理7.  $A \subset B$ ,  $A \in A$ および $C \subset \langle A \rangle - A$ とする時、 $C(\subset A)$ に対して次の条件1)～3)を考えると、条件1)から条件2)が、また条件3)から条件1)が成り立つ。

1)  $C$ は、 $C \ni A$ ,  $\langle C \rangle \cap C = \emptyset$ および $C \neq A$ を満たす、 $A$ における最大分解不能族である。

2)  $\langle A \rangle - A$ の元 $u$ および $C'$ が存在して、 $C'$ は $C' \supset C$ となる、 $A$ における $\{A\}$ の $(C' \cup \{u\}) - M I F$ である。

3)  $C \subseteq C' \subset \langle A \rangle - A$ なる $C'$ が存在して、 $C$ は $A$ における $\{A\}$ の $C' - M I F$ である。

定理8. 定理7を用いると、 $A$ における $\{A\}$ の $C' - M I F$ を求ることにより、定理8の条件1)を満たす、すべての $C$ を求めることができる。

これらの定理を用いて、次のような方法によって、 $\langle U, B, f \rangle$ に対するMSPDを求めることができる。

- (I)  $A (\subset B)$  のある元  $A$  に対して、 $A$  における  $\{A\}$  の  $\emptyset$ -極大分解不能族  $C$  をすべて求める。
- (II) (I) で求めた  $C$  の集合を用いて、すべての  $A$  の最大分解不能族  $D$  を求める。
- (III) (II) で求めた  $D$  の集合を用いて、すべての  $A$  の分解不能族を求める。
- (IV) (I) における  $A$  および  $A$  を変えることにより、すべての  $B$  の分解不能族を求める。
- (V) (IV) で求めた  $E$  の集合を用いて、分解不能族による  $B$  の分割の集合  $IP(B)$  を求める。
- (VI)  $IP(B)$  の各元に対し、条件 (C8), (C9) をチェックし、それらの条件を満たす元  $\beta$  の集合  $P(B, f)$  を求める。
- (VII) (VI) で求めた  $P(B, f)$  の元  $\beta$  に対し  $SP_f(A_\beta)$  を計算することにより、 $\langle U, B, f \rangle$  に対する  $MSPD$  を求める。

例 3.4. 図 3.4 は  $\langle U, B, f \rangle$  が与えられた時の  $MSPD$  の候補として考えられる  $B$  および  $\{U\}$  と、 $MSPD$  である。  $A_\theta$  の  $SP_f$  の値を比較したものである。

	B	U							B	MSPD の候補	
		u <sub>1</sub>	u <sub>2</sub>	u <sub>3</sub>	u <sub>4</sub>	u <sub>5</sub>	u <sub>6</sub>	u <sub>7</sub>		A <sub>θ</sub>	{U}
B	B <sub>1</sub>	1	1	1	1				80	160	1600
	B <sub>2</sub>	1	1	1		1			80		
	B <sub>3</sub>	1	1		1	1			80		
	B <sub>4</sub>	1	1			1	1		200	200	
	B <sub>5</sub>	1		1	1	1			80		
	B <sub>6</sub>		1	1	1	1			16		
	B <sub>7</sub>	1	1	1		1			40	80	80
	B <sub>8</sub>	1	1		1	1			40		
	B <sub>9</sub>	1		1	1	1			40		
	B <sub>10</sub>		1	1	1	1			40		
	B <sub>11</sub>	1	1	1			1		16		
	B <sub>12</sub>	1	1		1	1			16	32	32
	B <sub>13</sub>	1		1	1	1			16		
	B <sub>14</sub>		1	1	1	1			16		
f		10	2	2	2	2	5	2			
SP <sub>f</sub>									760	472	1600

図 3.4.

上に示した  $MSPD$  を求める方法 (I) ~ (VII) に従って、 $A_\theta$  を求める。

- 1)  $D^*(\{B_i\}, \emptyset) = \{B_i : i = 1, \dots, 10\} (= C_1 \text{ とおく})$  であり、 $C_1$  は分解不能族だから、 $B$  における  $\{B_i\}$  の  $\emptyset$ -MIF は  $C_1$  である。また  $\langle C_1 \rangle = \{u_i : i = 1, \dots, 6\}$  より、 $B$  における  $\{B_i\}$  の  $\{u_7\}$ -MIF でもある。
- 2)  $B$  における  $\{B_i\}$  の  $\{u_5\}$ -MIF は  $\{B_1\}$  である。 $B$  における  $\{B_i\}$  の  $\{u_6\}$ -MIF は、 $\{B_1, B_2, B_3, B_5, B_6\}$  ( $= C_2 \text{ とおく}$ ) である。 $\langle C_2 \rangle = \{u_i : i = 1, \dots, 5\}$  より  $B$  における

$\{B_1\}$  の  $\{u_6, u_7\}$  -MIF である。

- 3) 1), 2) より  $B$  における  $\{B_1\}$  の最大分解不能族は、 $\{B_1\}, \mathbb{C}_1, \mathbb{C}_2$  である。  
 4)  $B - \{B_1\}$  における  $\{B_2\}$  の  $\emptyset$ -MIF は、 $\mathbb{C}_1 - \{B_1\}$  である。 $\langle \mathbb{C}_1 - \{B_1\} \rangle = \{u_i : i = 1, \dots, 6\}$  より、 $B - \{B_1\}$  における  $\{B_2\}$  の  $\{u_7\}$  -MIF も、 $\mathbb{C}_1 - \{B_1\}$  である。  
 5)  $B - \{B_1\}$  における  $\{B_2\}$  の  $\{u_4\}$  -MIF および  $\{u_4, u_7\}$  -MIF は、 $\{B_2, B_4, B_8\}$  ( $= \mathbb{C}_3$  とおく) である。 $B - \{B_1\}$  における  $\{B_2\}$  の  $\{u_6\}$  -MIF および  $\{u_6, u_7\}$  -MIF は、 $\{B_2, B_3, B_5, B_6\}$  ( $= \mathbb{D}$  とおく) である。 $B - \{B_1\}$  における  $\{B_2\}$  の  $\{u_4, u_6\}$  -MIF は、 $\{B_2\}$  である。  
 6)  $\langle \mathbb{C}_1 - \{B_1\} \rangle = \langle \mathbb{C}_1 \rangle$  および  $\langle \mathbb{D} \rangle = \langle \mathbb{C}_2 \rangle$  だから、3), 4), 5) より、 $B$  における  $\{B_2\}$  の最大分解不能族は、 $\{B_2\}, \mathbb{C}_1, \mathbb{C}_2, \mathbb{C}_3$  である。  
 7) 以後同様の方法で、 $B$  における  $\{B_i\}$  ( $i = 1, \dots, 14$ ) のいづれかの最大分解不能族は、 $\{B_i\}$  ( $i = 1, \dots, 14$ ),  $\mathbb{C}_1 = \{B_i : i = 1, \dots, 10\}$ ,  $\mathbb{C}_2 = \{B_1, B_2, B_3, B_5, B_6\}$ ,  $\mathbb{C}_3 = \{B_2, B_4, B_8\}$ ,  $\mathbb{C}_4 = \{B_3, B_4, B_9\}$ ,  $\mathbb{C}_5 = \{B_i : i = 6, \dots, 10\}$ ,  $\mathbb{C}_6 = \{B_i : i = 6, 11, \dots, 14\}$  である。  
 8)  $\{\beta \in IP(B) : \beta \text{ は条件 (C8) を満たす}\} = \{\alpha_i : i = 1, \dots, 11\}$ ,  
 $\alpha_1 = \{\{B_i\} : i = 1, \dots, 14\}$ ,  $\alpha_2 = \{\mathbb{C}_1\} \cup \{\{B_i\} : i = 11, \dots, 14\}$ ,  
 $\alpha_3 = \{\mathbb{C}_1\} \cup \{\mathbb{C}_6 - \{B_6\}\}$ ,  $\alpha_4 = \{\mathbb{C}_2\} \cup \{\{B_i\} : i = 4, 6, \dots, 14\}$ ,  
 $\alpha_5 = \{\mathbb{C}_2, \mathbb{C}_5\} \cup \{\{B_i\} : i = 4, 11, \dots, 14\}$ ,  $\alpha_6 = \{\mathbb{C}_2, \mathbb{C}_5, \mathbb{C}_6 - \{B_6\}\} \cup \{\{B_4\}\}$ ,  
 $\alpha_7 = \{\mathbb{C}_3\} \cup \{\{B_i\} : i = 1, 3, 5, 6, 7, 9, \dots, 14\}$ ,  
 $\alpha_8 = \{\mathbb{C}_4\} \cup \{\{B_i\} : i = 1, 2, 5, \dots, 8, 11, \dots, 14\}$ ,  
 $\alpha_9 = \{\mathbb{C}_5\} \cup \{\{B_i\} : i = 1, \dots, 5, 11, \dots, 14\}$ ,  
 $\alpha_{10} = \{\mathbb{C}_5, \mathbb{C}_6 - \{B_6\}\} \cup \{\{B_i\} : i = 1, \dots, 5\}$ ,  
 $\alpha_{11} = \{\mathbb{C}_6\} \cup \{\{B_i\} : i = 1, \dots, 5, 7, \dots, 10\}$  となる。

$IP(B)$  の元であるが、条件 (C8) を満たさない  $\beta$  は多くある。例えば  $\mathbb{C}_2$  および  $\mathbb{C}_3$  を用いて  $IP(B)$  の元  $\beta$  を作ろうとすると、

- i)  $\beta \ni \mathbb{C}_2, \mathbb{C}_3 - \{B_2\}$  または  
 ii)  $\beta \ni \mathbb{C}_2 - \{B_2\}, \mathbb{C}_3$

のいづれかが成り立つ。i) の場合は  $\mathbb{C}_3 - \{B_2\}$  は分解可能となる。またii) の場合は、(C8) において  $\mathbb{C} = \mathbb{C}_3$ ,  $\mathbb{C}' = \mathbb{C}_2 - \{B_2\}$ ,  $A = \{B_2\}$  とおくと、 $A \subset \mathbb{C}$ かつ  $\langle A \rangle \subset \langle \mathbb{C}' \rangle$  が成り立つが、 $\mathbb{C} - A$  は分解可能になる。他の場合も同様にチェックできる。

- 9)  $SP_f(\langle \mathbb{C}_1 \rangle) (= 800) > SP_f(\mathbb{C}_1) (= 696)$ ,  
 $SP_f(\langle \mathbb{C}_2 \rangle) (= 160) < SP_f(\mathbb{C}_2) (= 320)$ ,  
 $SP_f(\langle \mathbb{C}_3 \rangle) (= 400) > SP_f(\mathbb{C}_3) (= 320)$ , ( $i = 3, 4$ )  
 $SP_f(\langle \mathbb{C}_5 \rangle) (= 80) < SP_f(\mathbb{C}_5) (= 176)$ ,  
 $SP_f(\langle \mathbb{C}_6 \rangle) (= 32) < SP_f(\mathbb{C}_6) (= 80)$ ,  
 $SP_f(\langle \mathbb{C}_6 - \{B_6\} \rangle) (= 32) < SP_f(\mathbb{C}_6 - \{B_6\}) (= 64)$

だから  $P(B, f) = \{\alpha_i : i = 1, 4, 5, 6, 9, 10, 11\}$  となる。

- 10)  $SP_f(A_{\alpha_1}) = 760$ ,  $SP_f(A_{\alpha_4}) = 600$ ,  $SP_f(A_{\alpha_5}) = 504$ ,  
 $SP_f(A_{\alpha_6}) = 456$ ,  $SP_f(A_{\alpha_9}) = 664$ ,  $SP_f(A_{\alpha_{10}}) = 616$ ,  
 $SP_f(A_{\alpha_{11}}) = 712$  より、 $A_{\alpha_6}$  ( $= A_0$ ) が MSPD となる。

#### 4. まとめ

我々は、統計データベースの設計問題の定式化と解を求めるための手法を与えた。しかし、これで統計データベースの設計の問題が完全に解決したわけではない。次に残された問題について述べる。

1)  $S_{G,x}$  の定義でわかるように、ここでの操作は射影と集約のみを考えている。結合によってできる要約表は、元々素データ  $R$  より射影と集約で表現できるため、ここでは取り扱う必要がないと考えた。今後選択を含めた問題の定式化をする必要がある。

2) 定理 1において、 $S_0$  の任意の元  $S_{G,x}$  に対し  $G = G_0$  を仮定したが、一般的にはこのような仮定はできない。従ってこの仮定が成立しない場合の問題の解法を求める必要がある。

3) 要約表の計算機での格納容量としてレコード数を考えた。これは分類項目の値は索引として各項目ごとに異なる値のみを記憶できること、および 2) における仮定があるため、格納容量はレコード数の定数倍と考えても不自然でないためであった。従って 2) の仮定をはずした場合には別の値を考える必要がある。

4)  $|B|$  や  $|U|$  の値が大きくなつた時、MSPD を求める方法は効率的かといった心配がある。統計データベースでは  $|U| \leq 20$ 、 $|B| \leq 200$  を満たす場合がほとんどであるため、この範囲内では我々の求めた方法は、MSPD を求めるために効果的であると考えている。今後この方法の改良および評価を行う必要がある。

[謝辞] 本稿作成において、多くの助言をいただいた、広島大学の菅原正博教授および池田秀人助教授に深く謝意を表します。

- [参考文献]
- [1] Y. Kobayashi and H. Ikeda, Statistical data model and semantics, Proc. of World Conference on Information Processing/Communication (1989), 438-443.
  - [2] Y. Kobayashi, A mathematical study on statistical database designs, to appear in Hiroshima Mathematical Journal.