

# 実世界データベース(2)

2J-3

-生成・概念・表記法-

池田 一夫、倉科 周介、大橋 誠

東京都立衛生研究所

## 1. はじめに

国や地方自治体では、社会・経済の現状把握や政策立案の基礎資料として、各種統計データを収集している。1988年には、統計データ収録管理用データベース(DB)の設計開発手法が佐藤<sup>1)</sup>により提案されている。

統計データのような一次知識(primary knowledge: 外界に存在する対象を直接観測することによって得られる知識)は、調査対象および調査目的・調査結果の編集法や記録法などによって規定される。特に時系列データのように広範な展開規模を持つデータの場合、個別調査はその時点の断片的な情報を与えるにとどまる。したがって、各調査データを総合的に検討することにより得られる知識は、現状の方法に依存する限り生産することができない。

本論文では、この制約を克服するための二次知識(secondary knowledge)の生成法、その結果生成する実世界DB(real world database)およびその表記法について論じる。

## 2. 二次知識の生成法と実世界DBの概念

### 2-1. 一次知識の蒐集

対象とする事物に関する一次知識を網羅的に蒐集する。蒐集した一次知識(主としてデータ)は、記録の原型を保ったまま蓄積する(primary data file)。この段階では、内容の修正などは原則として一切行わない。

### 2-2. 仮想空間の設定

実世界の時空間と同型の座標系を設定する。実世界の事物はすべてこの座標系で固有の位置を占める(real world framework)。さらにこの座標系を基準にして、事物の状態を特徴づける座標系を組み込む。人口現象や疾病現象を対象とする場合は、生物学的座標軸(性別、年代、年齢、疾病類別など)、社会的座標軸(職業、組織、その他各種の社会的条件)等をデータの特性に応じて適切に付加する。

### 2-3. データの再配置

primary fileのデータの正誤を照合調整して、運用のためのデータファイル(revised file)を作成する。これから個別データを抽出して real world frameworkのしかるべき位置に配置する(real world data file)。この段階では、データの記述単位は一次知識のそれと一致する。

### 2-4. 実世界の復元

real world data fileのデータ構造をもとに、個々のデータを分解・再配置し、データの欠損部分は近傍データとの連続性およびデータ全体の中での整合性の両面に配慮しつつ推計補填して、過去の実世界の全体像を復元する(real world database)。復元に用いるデータは一次データではなく、それに基づいた人工的なもの(artificial data)である。

### 2-5. 実世界の観測

復元された実世界データを基準に、一次知識の生成の場合と同様な観測を行う。すなわち、人工的に構築した世界から実世界を間接観測する。これを二次観測、または知識観測と呼ぶことにする。

### 3. 表記法

実世界DBが取り扱うのは統計データである。統計データの基礎となるのは個々のデータそれ自身である。そこで、個々のデータの表記法を基に各種概念の表記法を示す。

n次元空間の各座標軸の名前を $A_1, A_2, \dots, A_n$ とし、各座標軸の任意の要素を $a_1, a_2, \dots, a_n$ と表記する。上位カテゴリと下位カテゴリ（都道府県と市町村など）を区別するために、座標軸および要素にレベルの概念を取入れ、それぞれ $A_n(L_n), a_n(L_n)$ と表記する。カッコ内の値はレベルを示し、標準状態では $L_n$ の値は0である。表記を簡略化するために、誤解が生じない限り(0)の表記は省略できるものとする。レベルが下がるほど $L_n$ の値は小さくなり。標準状態より下のレベルでは負の値を持つものとする。

個々のデータはd（小文字）、データの集合はD（大文字）で表記する。

#### 3-1. 一般型

n次元空間の任意の位置における任意のデータの値を

$d \quad \mathbf{a} \quad \mathbf{a} : \text{ベクトル} [a_1 \ a_2 \ \dots \ a_n]$  と表記する。

誤解を生じない限り  $\mathbf{a}$  を  $a_1 a_2 \dots a_n$  と表記してもよい。

#### 3-2. primary data file のデータおよび座標軸

一次知識であることを表示するための記号 0 を導入し、

データ： $d \quad \mathbf{a} \quad \mathbf{a}^0$       座標軸： $A_n \quad A_n^0$       と表記する。

#### 3-3. real world data file のデータ

primary data を加工した人工的なものであることを明記する必要がある。real world data を使用する際には、そのデータがどのような過程を経て生成されたものなのかを式で明示する。形式は次のとおりとする。

座標軸  $A_n \leftarrow A_n^0(L_n)$   
 $[A_1 \ A_2 \ \dots \ A_n] \leftarrow [A_1^0(L_1) \ \dots \ A_n^0(L_n)]$

データは、一般形式で表記する。

#### 3-4. イメージデータ

特定の座標軸 $A_i, A_k, \dots, A_m$ 上の特定の値 $a_i, a_k, \dots, a_m$ におけるデータのプロフィールを次のように表記する。

$D \mid \mathbf{a} \quad \mathbf{a} : \text{ベクトル} [0 \ 0 \ \dots \ a_i \ \dots \ a_k \ \dots \ a_m \ \dots \ 0 \ 0]$

誤解を生じない限り  $\mathbf{a}$  を  $0 \ 0 \ \dots \ a_i \ \dots \ a_k \ \dots \ a_m \ \dots \ 0 \ 0$  または単に  $a_i a_k \dots a_m$  と表記してもよい。

### 4. 表記例

時刻 T  $[t_1 \ t_2 \ \dots \ t_n]$       世代 G  $[g_1 \ g_2 \ \dots \ g_n]$   
 空間 S  $[s_1 \ s_2 \ \dots \ s_n]$       要素 E  $[e_1 \ e_2 \ \dots \ e_n]$

として、各種の表記例を示す（[]内は各々の要素）。

●一般的なデータ： $d_{tsge}$       (t, s, g, e は要素の代表)

●特定時刻・特定空間・特定要素の世代構成： $D \mid tse$

●市区町村データを利用して都道府県ごとの合計を計算する。

$$d_{ts(1)ge} = \sum_s d_{tsge}$$

●さらに世代別合計をとる。

$$d_{ts(1)g(1)e} = \sum_g d_{ts(1)ge} = \sum_g \sum_s d_{tsge}$$

#### [参考文献]

1) 佐藤英人著、統計データベースの設計と開発、オーム社、1988