

## A mathematical study on statistical database designs

Yasuyuki KOBAYASHI

(Received September 4, 1989)

### 1. Introduction

Let be given a finite set  $U$  and non-negative integers  $f(x)$  for all  $x \in U$ . Then, by taking the sum of products of them, we have an integer

$$(1) \quad \text{SP}_f(\mathcal{A}) = \sum_{A \in \mathcal{A}} \prod_{x \in A} f(x)$$

for each subfamily  $\mathcal{A} \subset 2^U - \{\emptyset\}$ , especially, for any covering  $\mathcal{A}$  of  $U$ ; and we can consider the following

**PROBLEM.** For given  $U, f$  as above and a covering  $\mathcal{B}$  of  $U$ , find effectively  $\mathcal{A}$  in such coverings  $\mathcal{A}'$  that  $\mathcal{B}$  is a refinement of  $\mathcal{A}'$  so that the function  $\text{SP}_f$  in (1) takes the minimum value at  $\mathcal{A}$  among such coverings  $\mathcal{A}'$  (see Definition 2.3).

We call  $\mathcal{A}$  in this problem an MSPD for  $\langle U, \mathcal{B}, f \rangle$  simply. Of course, an MSPD exists and any MSPD can be found by calculating  $\text{SP}_f(\mathcal{A}')$  for all finitely many such  $\mathcal{A}'$ ; but the number of  $\mathcal{A}'$  may increase rapidly as  $|U|$  increases. ( $|X|$  denotes the number of elements in a finite set  $X$ .)

Thus, the purpose of this paper is to establish an effective method of finding an MSPD of special type, which is applicable even when  $|U|$  may be large.

Our motivation is in the problem on statistical database designs stated in § 5. (For databases, cf. Codd [3–5] and Smith-Smith [23], and for statistical databases, cf. Shoshani [22] and several papers in the reference.)

Let  $R$  be a given collection of statistical records, that is, a finite subset of the product  $D = \prod_{i=1}^N D_i$  of domains  $D_i$  of  $i$ -th field. Then, an aggregation function  $S$  can be specified by the category fields  $X(S)$ , the summary fields  $Y(S)$  and the summarizing operators  $g_j$  over  $D_j$  given for each summary field  $j$  in  $Y(S)$ ; and  $S$  gives us the summary table  $S(R)$  corresponding to  $X(S)$ ,  $Y(S)$  and  $g_j$ 's. Moreover, for any finite set  $\mathcal{S}$  of aggregation functions, we have

$$(2) \quad \text{NRec}(\mathcal{S}) = \sum_{S \in \mathcal{S}} |S(R)|,$$

the total number of records of  $\{S(R): S \in \mathcal{S}\}$ . Thus we have the following

**PROBLEM.** Let  $R$  be a given collection of statistical records. Then, for a finite set of summary tables  $\{S_0(R): S_0 \in \mathcal{S}_0\}$  to be derived from the database,