

Purdue University

Purdue e-Pubs

Department of Computer Science Technical
Reports

Department of Computer Science

1976

Linear Queries in Statistical Data Bases

M. D. Swartz

D. E. Denning

P. J. Denning

Report Number:

77-216

Swartz, M. D.; Denning, D. E.; and Denning, P. J., "Linear Queries in Statistical Data Bases" (1976).
Department of Computer Science Technical Reports. Paper 156.
<https://docs.lib.purdue.edu/cstech/156>

This document has been made available through Purdue e-Pubs, a service of the Purdue University Libraries.
Please contact epubs@purdue.edu for additional information.

LINEAR QUERIES IN STATISTICAL DATA BASES*

M. D. Schwartz
D. E. Denning
P. J. Denning

Purdue University
Department Computer Science
W. Lafayette, IN 47907

November 1976
Revised February 1978

CSD-TR-216

Abstract. A data base is compromised if a user can determine the data elements associated with keys which he did not know previously. If it is possible, compromise can be achieved by posing a finite set of queries over sets of data elements and employing initial information to solve the resulting system of equations. Assuming the allowable queries are linear, that is, weighted sums of data elements, we show how compromise can be achieved and we characterize the maximal initial information permitted of a user in a secure system. When compromise is possible, the initial information and the number of queries required to achieve it is surprisingly small.

* Work reported herein was supported in part by NSF Grant MCS75-21100 and by Purdue David Ross Grant 7994-56-13985.

1. INTRODUCTION

Statistical data bases contain sensitive information about individuals. Their purpose is to provide summary statistics about groups of people, while permitting only authorized access to the records of any one individual. But this objective is difficult to meet, for seemingly innocuous summaries contain vestiges of the original data. By correlating enough summaries, an intruder may infer confidential data. We understand much more about controlling access to individual records than we do about controlling inference.

A query is a function applied to a given subset of the records in the data base; this subset is called the query set. Given the form of allowable queries and an assumption about the information known initially by enquirers, we wish to specify whether or not the data base is secure. A proof that a data base can be compromised amounts to exhibiting a method for constructing a sequence of queries whose responses imply the value associated with a given, arbitrary key. A proof that a data base is secure is much more difficult, for it requires us to demonstrate that no such sequence exists.

Inference has long been a concern of census statisticians, who have developed sophisticated methods for choosing samples for analysis. The principle is to apply queries against a random subfile; deducing an element of the subfile is of little help because the relation between the subfile and the full population is unknown. Unfortunately, many contemporary applications center on small or medium data bases where random sampling cannot be used.

Much existing literature on inference applies to counting queries, which return the number of individuals satisfying given characteristics -- e.g., "Female astronomers over 40 years of age." [Haq74, Haq75, HoM70, Sch75]. These studies have shown that the danger of compromise is greatest when the system responds for small query sets. Some protection is afforded when the system refuses answers for small counts [Chi77, FeP74, Sch76]. However, this protection is far from complete [DDS77, Scw77].

Another class of easily compromised queries is selection queries. Such a query returns a specific value from the query set, such as the largest, the smallest, or the median. These queries are easily compromised even if the data base system somehow enforces a restriction that no two query sets overlap by more than one record [Dav76, DJL76, Rei77, DDL76].

The subject of this paper is the class of linear queries. Each of these queries returns a weighted sum of elements in the query set. Our results extend those of Dobkin, Jones, and Lipton [DJL76], whose queries computed simple (unweighted) sums of elements in the query sets. One aspect of the results is striking: a user who knows as little as the value associated with a single key can often achieve a full compromise of the entire data base -- with a query sequence whose length is not much longer than the number of stored data elements.

As a running example, we will use the data base displayed in Table 1, referred to hereafter as the "runner's data base." It contains facts about foot racers. Users of the data base might be allowed to discover general properties such as average running paces; however, they are not supposed to be able to deduce facts about individual runners from these queries.

Table 1. Runner's Data Base.

KEYS i	j	1 Max VOX	2 Train pace sec/mi	3 Prior 8 weeks		5 Fastest 1 mile sec.	6 Fastest 10 miles sec.	7 Year of Birth
				Total Miles	Longest Run			
1	Smith	68	380	680	22	260	3183	1948
2	Jones	61	405	530	18	287	3520	1950
3	Burns	56	440	460	20	316	3818	1940
4	Cohen	48	485	410	18	368	4447	1930
5	Cook	49	470	375	20	360	4394	1962
6	Bloom	54	440	430	20	328	3980	1948
7	King	53	440	405	20	334	4072	1943
8	Frank	72	370	705	23	251	2997	1947

"Max VOX" is the maximum volume of oxygen processable
in milliliters per kilogram body weight per minute.

MODEL OF A DATA BASE

A statistical data base contains a set of records about individuals. Each record contains an identifier field, one or more category fields, and one or more data fields. Some data fields may also be category fields.

The identifier field contains a unique identifier for the record. Reading and writing in single records via identifiers is controlled by an access mechanism. A user presents the symbolic name of a record, the system looks up the identifier corresponding to that name and checks authorization.

A data element is a value stored in a given data field of a given record. The size of the data base is the total number, N , of data elements in it. In a partitioned data base, the data elements are divided into mutually exclusive, collectively exhaustive classes. The j th class, denoted Y_j , contains N_j data elements, and $N = \sum_j N_j$. The columns of the runner's data base are examples of classes. A key $p = (i,j)$ identifies a data element, which is denoted either x_p or x_{ij} .

A query set is a collection of data elements identified by a list of keys or by a characteristic formula. (A characteristic formula is a logical formula which is matched against the values in category fields.) A query computes some value for the query set. To prevent trivial compromises a query program will not respond if the query set contains fewer than k data elements, where $k \geq 2$ is a parameter of the system.

EXAMPLE. The runner's data base contains 8 records and 7 data fields, a total of $N = 56$ data elements. The keys are pairs (i,j) where i is a row index and j is a column index. The user would specify symbolic name pairs which would be translated to keys internally; for example, (Smith, Max VOX) translates to $(1,1)$ and (Cohen, Year of Birth) to $(4,7)$. The corresponding data elements are $x_{11} = 68$ and $x_{47} = 1930$. A possible key-specified query is "What is the average train-pace for (Smith, Jones, Frank)?" The same query expressed with a characteristic formula is "What is the average train-pace among runners whose Max VOX is greater than 60?"

This paper analyzes threats to statistical data base security for key-specified query sets. Security for characteristic-specified query sets is treated in other papers [DDS77, Scw77]. A summary of important results for both types of query sets is given in [Den78].

Compromising the Data Base

A user has identified a key (i,j) as soon as he knows the value of its data element x_{ij} . A user usually has some initial information about the data base: he may know operational details of some query programs, or the values of some data elements. Relative to an assumption about initial information, a data base is insecure if it is possible to identify keys not known initially. Identifying previously unknown keys is a compromise,

and identifying all the keys is a full compromise. If the data base is insecure there will exist a finite sequence of queries q_1, \dots, q_m , query sets X_1, \dots, X_m , and responses v_1, \dots, v_m for which it is possible to solve the system of equations

$$\begin{aligned}
 v_1 &= q_1(X_1) \\
 v_2 &= q_2(X_2) \\
 (1) \quad &\dots \\
 v_m &= q_m(X_m)
 \end{aligned}$$

for some unknown key.

We are interested in characterizing the amount of initial information that can be tolerated in a secure system, and the amount of work of compromise an insecure system.

EXAMPLE. A regression analysis has yielded a formula to predict the time a runner will take to run a marathon [FoD75]:

$$(2) \quad q(i) = \sum_{j=1}^4 a_j x_{ij} + 319 \text{ seconds.}$$

Suppose that this query is available in the data base and that the user knows the coefficients a_j (e.g., from reading the article).

If that user also knows any three of x_{i1} , x_{i2} , x_{i3} , and x_{i4} , he can use the system's response to $q(i)$ to solve Equation 2 for the fourth.

Linear Queries

To pose a query, a user specifies a list of distinct keys p_1, \dots, p_k ; the system determines the query set $\{x_{p_1}, \dots, x_{p_k}\}$ and computes the value of the query function for it. To keep the notation simple, we will show the query set as argument to the query, rather than the list of keys.

Thus $q(z_1, \dots, z_k)$ denotes a query over a list of distinct keys that specify the data elements z_1, \dots, z_k . A linear query has the form

$$(3) \quad q(z_1, \dots, z_k) = \sum_{j=1}^k a_j z_j$$

for some fixed $k > 2$ and fixed "query weights" a_j . This form of query has been called a "weighted sum query" [SDD77].

The assumption that the query set size is fixed at k elements for all queries is a degenerate case of the restriction that each query set contains at least k elements. When query set size can vary, as for characteristic-specified queries, compromise can be even simpler [DDS77]. Our objective here is to show how easy compromise can be even when the intruder cannot exploit variation in query set size.

A sum query is a linear query with unit weights (all $a_j = 1$). Such queries are used to compute averages. Dobkin, Jones, and Lipton [DJL76] developed upper and lower bounds on the number of sum queries sufficient to compromise under the restriction that no two query sets can overlap by more than r elements. For N data elements with $r = 1$ and no initial information, compromise is possible within $2k-1$ queries if $N \geq N_2 = k^2 - k + 1$; however, compromise is impossible when $N < N_1 = (k^2 + k)/2$. Whether or not compromise is possible when $N_1 \leq N < N_2$ was not resolved. Davida et al tightened these bounds [Dav76].

Kam and Ullman studied the security of a class of key-matching sum queries [KaU76]. Each key is a bit-string of b bits, and there are 2^b data elements. To make a query, the user specifies the values of some set of a bits, where $a \leq b$; the query program returns the sum of the data elements whose keys match in the given a positions. The data base is secure whenever users have no initial information, if $a < b$, and if the ranges of values of the data elements are unknown. Otherwise compromise may be possible.

A partitioned query is a linear query in a partitioned data base with the restriction that the j th key specifies a data element from the j th class; there are k classes in all. Equation 2 is an example with four classes.

Our results for linear queries and partitioned queries are outlined in the next sections. The proofs are given in the appendices.

3 SECURITY OF LINEAR QUERIES

We consider queries of the form of Equation 3 applied to a data base of N elements. All queries use the same value of k and the same weights a_j .

We observe first that, knowing one weight a_j and one data element x , we can compromise the entire data base. Let q_1 and q_2 be the responses to the two queries

$$\begin{aligned} q_1 &= q(z_1, \dots, z_{j-1}, x, z_{j+1}, \dots, z_k) \\ q_2 &= q(z_1, \dots, z_{j-1}, y, z_{j+1}, \dots, z_k) \end{aligned}$$

which differ only in their use of data elements x and y in the j th position. Equation 3 shows that

$$(4) \quad (q_1 - q_2) = a_j(x - y),$$

which can be solved for y . Now the entire data base is vulnerable: we pose q_1 once and a new q_2 for each of the $N-1$ unknown data elements, effecting a full compromise with N queries and $N-1$ applications of Equation 4. As soon as two data elements, x and y , are known, we can also use Equation 4 to solve for any unknown weight.

EXAMPLE. Suppose that the runner's data base implements the query

$$q(z_1, z_2, z_3) = .2z_1 + .5z_2 + .3z_3$$

A user knows that the weight of the first key is .2 and that Smith's Max VOX is 68; thus the initial information is

$$a_1 = .2$$

$$x_{11} = 68$$

To determine Jones's Max VOX, the user proceeds as follows. He poses the two queries $q(x_{11}, x_{12}, x_{13})$ and $q(x_{21}, x_{12}, x_{13})$, to which the system will respond, respectively,

$$q_1 = 407.6$$

$$q_2 = 406.2$$

Equation 4 can be used to solve for x_{21} :

$$\begin{aligned} x_{21} &= x_{11} - (q_1 - q_2)/a_1 \\ &= 68 - 1.4/.2 \\ &= 68 - 7 \\ &= 61 \end{aligned}$$

which is Smith's MAX VOX.

This reasoning suggests that security, if it is possible at all, must require keeping all the weights secret. But a more complex argument, given in Appendix 1, shows that we can pose a set of queries which yield a set of equations in the unknowns a_j . If these equations are linearly independent, we can solve for the a_j and then use the known value of one data element to compromise the rest of the data base. In other words, knowing only one data element and none of the weights often suffices for a compromise. The exact statement of the result is:

Compromising Weighted Sums. Suppose that x_p is known for some key p . With no more than $k(k+1)$ queries it may be possible to determine x_{p_1}, \dots, x_{p_k} for k additional keys p_1, \dots, p_k . The remaining $N-k-1$ data elements can be deduced with an additional $N-k-1$ queries using the method of Equation 4. Thus full compromise is possible within $N+k^2-1$ queries.

The technique will be illustrated in an example below.

These arguments lead to the conclusion that security depends on denying the inquirer any initial information about the data base -- for example, by restricting his access to records about which he knows nothing. In Appendix 2 we prove that such a restriction would work by showing the impossibility of solving any system of equations for any weight or any data element:

Secure Weighted Sums. If the inquirer knows no data element, the data base is secure under his queries.

Unfortunately there is no way to be absolutely sure that a user knows nothing about the records to which his access has been restricted. Even if such an idea were enforceable, the presence of other types of queries can invalidate this result. The unweighted-sum compromise of Dobkin et al, for example, requires no initial knowledge [DJL76].

The next subsection illustrates the method of compromising a data base using weighted sum queries when one data element is known.

Example.

Let $y_i = x_{i1}$ denote the Max VOX for runner i . Let $k = 2$ and suppose $y_1 = 68$ is known. The data base implements the query

$$q(z_1, z_2) = 2z_1 + 3z_2$$

but the weights $(a_1, a_2) = (2, 3)$ are unknown. Form three sets of two queries:

$$\begin{aligned}
 q_{10} &= q(y_2, y_3) = 290 & q_{11} &= q(y_3, y_2) = 295 \\
 q_{20} &= q(y_1, y_3) = 304 & q_{21} &= q(y_3, y_1) = 316 \\
 q_{30} &= q(y_1, y_2) = 319 & q_{32} &= q(y_2, y_1) = 326
 \end{aligned}$$

the numbers show the system's responses. Let $A = a_1 + a_2$ and observe that

$$\begin{aligned}
 q_{10} + q_{11} &= A(y_2 + y_3) = 585 = Q_1 \\
 q_{20} + q_{21} &= A(y_1 + y_3) = 620 = Q_2 \\
 q_{31} + q_{32} &= A(y_1 + y_2) = 645 = Q_3
 \end{aligned}$$

whereupon

$$A = Q_1 / (y_2 + y_3) = Q_2 / (y_1 + y_3) = Q_3 / (y_1 + y_2)$$

This can be manipulated into two equations in unknowns y_1 and y_2 :

$$\begin{aligned}
 -Q_2 y_2 + (Q_1 - Q_2) y_3 &= -Q_1 y_1 \\
 Q_2 y_2 - Q_3 y_3 &= (Q_3 - Q_2) y_1
 \end{aligned}$$

For the given numbers,

$$\begin{aligned}
 -620 y_2 - 35 y_3 &= -(585) \quad (68) \\
 620 y_2 - 645 y_3 &= (25) \quad (68)
 \end{aligned}$$

It is easily verified that $y_2 = 61$ and $y_3 = 56$ is the solution. We may solve for a_1 and a_2 using the equations

$$\begin{aligned}
 a_1 y_2 + a_2 y_3 &= q_{10} \\
 a_1 y_1 + a_2 y_3 &= q_{20}
 \end{aligned}$$

To find any other y_i we solve $q(y_1, y_i) = 136 + 3y_i$ once the system's response is known.

4. PARTITIONED QUERIES

We consider weighted-sum queries with the restriction that each weight a_j is associated with a separate class Y_j of data elements, and that the j^{th} key in a query's specification list must refer to the class Y_j . We denote the number of data elements in Y_j by N_j , so that $N = \sum_j N_j$. The columns of the runner's data base are examples of classes.

We observe first that the method of Equation 4 can be used to compromise all of class Y_j if the weight a_j and some element x of that class both are known. In all, N_j queries are required for this. If we know all the weights and one data element in each class, we can compromise the entire data base with N queries.

Suppose that we know all k of the weights and data elements from only $k-1$ classes. We can determine an element from the unknown class (say Y_j) with one query using the known elements from the other classes (say $z_1, \dots, z_{j-1}, z_{j+1}, \dots, z_k$):

$$(5) \quad q(z_1, \dots, z_k) = a_1 z_1 + \dots + a_k z_k.$$

After the system responds with the value of this query, z_j is the only unknown quantity in this equation. It follows that full compromise is possible with $N+1$ queries when all weights are known and data elements from all but one class are known.

Similarly, if we know one element from each class, and all but one of the weights, we can apply Equation 5 using the known z_1, \dots, z_k ; the missing weight is the only unknown quantity. Full compromise is again possible with $N+1$ queries.

The method illustrated in the previous section was able to effect a compromise when none of the weights was known and the entire data base was a single class. This method depends on permuting the same keys among different positions in the query specification lists. Partitioning the data bases prevents these permutations and makes compromise more difficult. In Appendix 3 we prove that

If one knows all the weights, one data element from each of $k-2$ classes, and no data element from the other two classes, then full compromise of the $k-2$ known classes is possible. The two unknown classes are secure.

The proof shows that no system of equations constructed contains enough information to solve for data elements in the unknown classes. The proof of this result shows further that every unknown class is secure provided that at least two classes are unknown initially.

These results show that keeping the weights secret is essential to security especially if users may know data elements from some of the classes. However, if other types of queries are possible in the same data base, compromise may be possible without prior knowledge. For example, DeMillo et al have shown that $O(k^2)$ selection queries (median, maximum, minimum, etc.) can determine an element from an unknown class [DDL76]; such queries can provide the initial information needed to compromise with partitioned queries. (See also [Scw77].)

6. CONCLUSIONS

Our results have shown that the maximum tolerable information to prevent compromise using weighted-sum or partitioned-weighted-sum queries is quite low. Security under weighted-sum queries can be assured only if users have no prior knowledge of any data element; the weights must be kept secret. Security in partitioned data bases can be assured only if at least two of the weights or at least two of the classes are kept secret. It is not always possible to keep weights secret, and it may be impossible to enforce the restriction that the user knows no data element.

When the inquirer has sufficient initial information for a compromise, his task is easy. At worst, he needs to pose one query for each unknown data element in the (class of the) data base which he is compromising, plus possibly a small additional but fixed number of queries. If the data base supports other types of queries such as selection queries, the user may be able to deduce unknown data elements with small cost, before using algorithms of this paper for a compromise.

Compromise of statistical data bases is usually possible, and it is cheap.

APPENDIX 1: Compromising Linear Queries

We wish to show that knowing one data element, say z_1 , may be sufficient to compromise a data base under linear queries. The compromise proceeds in three parts. First, $k(k+1)$ queries are used to derive $k+1$ equations whose solution yields k additional data elements z_2, \dots, z_{k+1} . Then the $k+1$ known data elements are used to derive k more equations whose solution is the weights a_1, \dots, a_k . Finally, the method of Equation 4 is used to deduce the remaining $N-(k+1)$ unknown data elements. Full compromise is achievable with $N-(k+1)+k(k+1) = N+k^2-1$ queries in all.

Let $\underline{Z} = (z_1, \dots, z_{k+1})$ be a list of data elements for which z_1 is known. Let \underline{Z}_i be \underline{Z} with z_i deleted, and let \underline{Z}_{ij} denote the j th one-step cyclic permutation of \underline{Z}_i . (Note that $\underline{Z}_{i0} = \underline{Z}_{ik}$.) Let q_{ij} denote the response to the query $q(\underline{Z}_{ij})$, where $i = 1, \dots, k+1$ and $j = 0, \dots, k-1$. Note that $k(k+1)$ queries are needed to determine all the q_{ij} .

Each element of \underline{Z}_i is permuted into each of the k positions by the queries $q(\underline{Z}_{i0}), \dots, q(\underline{Z}_{i,k-1})$. If $S = z_1 + z_2 + \dots + z_{k+1}$ then the sum of elements in \underline{Z}_i is $S - z_i$ and the sum of these permuted queries is

$$Q_i = \sum_{j=0}^{k-1} q_{ij} = \sum_{j=0}^{k-1} a_j (S - z_i) = A(S - z_i)$$

where $A = a_1 + \dots + a_k$. It follows that $Q_i / (S - z_i) = A = Q_{i+1} / (S - z_{i+1})$. This gives a set of k equations in the k unknowns z_2, \dots, z_{k+1} :

$$S(Q_i - Q_{i+1}) + Q_{i+1}z_i - Q_i z_{i+1} = 0 \quad (i = 1, \dots, k)$$

If these equations are linearly independent, their solution will yield the values of z_2, \dots, z_{k+1} . (If these equations are not linearly independent, nothing can be concluded about the security of the data base, because another choice of z_2, \dots, z_{k+1} might produce a solvable set of equations.) We can use the known values of z_1, \dots, z_k to invert the equations corresponding to the queries $q(z_{i0}), \dots, q(z_{ik})$ and find the unknown coefficients:

$$\begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_k \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_k \end{bmatrix} = \begin{bmatrix} q_{10} \\ q_{20} \\ \vdots \\ q_{k0} \end{bmatrix}$$

Now the method of Equation 4 can be applied to find the remaining data elements.

APPENDIX 2: Secure Linear Queries

We suppose that the user has no initial knowledge about any data element. Our objective is proving that it is impossible for him to find enough linearly independent equations to solve for any unknown data element or query weight.

Consider any m queries q_1, \dots, q_m , among which appear the data elements $Z = \{z_1, \dots, z_s\}$. We assume $m \geq k+s$ since we can always pose more queries without impairing our ability to compromise. We will apply the Implicit Function Theorem (IFT) [Gar68], which states that a set of functions q_1, \dots, q_m over variables x_1, \dots, x_n are mutually independent if the $\det(M) \neq 0$, where

$$M = \begin{bmatrix} \frac{\partial q_1}{\partial x_1} & \dots & \frac{\partial q_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial q_m}{\partial x_1} & \dots & \frac{\partial q_m}{\partial x_n} \end{bmatrix}$$

For our problem, the queries q_1, \dots, q_m depend on the $n = k+s$ quantities a_1, \dots, a_k and z_1, \dots, z_s . The matrix is

$$M = \left[\begin{array}{ccc|ccc} z_{11} & \dots & z_{1k} & b_{11} & \dots & b_{1s} \\ & & & & & \\ & & & & & \\ z_{m1} & \dots & z_{mk} & b_{m1} & \dots & b_{ms} \end{array} \right]$$

where z_{ij} is the member of Z appearing in the j th position of the i th query, and

$$b_{it} = \begin{cases} a_j, & \text{if } z_{ij} = z_t \\ 0, & \text{otherwise} \end{cases}$$

Let $r = \text{rank}(M)$. We observe that $r < k+s$ because column $k+p$ ($p = 1, \dots, s$) can be written as a linear combination of the others:

$$\sum_{j=1}^k a_j z_{ij} - \sum_{\substack{t=1 \\ t \neq p}}^s z_t b_{it} = b_{ip} z_p \quad (i = 1, \dots, m)$$

Note that $z_{ij} = z_t$ if and only if $b_{it} = a_j$, which implies that the nonzero terms in the second sum cancel all but $b_{ip} z_p$ in the first. Since $r < k+s$, $\det(M) = 0$, which implies that some of q_1, \dots, q_m are redundant.

Consider an $r \times r$ submatrix N of M formed by taking any r linearly independent rows, and any r linearly independent columns except column $k+p$. Renumber the queries so that q_1, \dots, q_r correspond to the rows of N . By construction $\det(N) \neq 0$, but any determinant of $r+1$ rows and columns of M must be zero. The IFT shows that q_1, \dots, q_r are mutually independent and that each query among q_{r+1}, \dots, q_{k+s} can be expressed in terms of q_1, \dots, q_r [Gar68].

We will show that x_p is not a function of any of the queries q_1, \dots, q_r . To do this we treat x_p as a new function of the same quantities (a_j or z_t) on which q_1, \dots, q_r depend according to the matrix N ; we add the appropriate $r+1$ st row and column of N :

$$N' = \left[\begin{array}{cccc|c} & & & & \frac{\partial q_1}{\partial x_p} \\ & & & & \vdots \\ & & N & & \frac{\partial q_r}{\partial x_p} \\ \hline & & & & \\ \hline 0 & 0 & \dots & 0 & 1 \end{array} \right]$$

It is not hard to see that $\det(N') = (-1)^r \det(N) \neq 0$, which implies that x_p is not a function of any quantity on which q_1, \dots, q_r depend and, therefore, that x_p is not a function of any of q_1, \dots, q_r .

APPENDIX 3: Secure Partitioned Queries

We wish to show that if all weights are known but no element from at least two classes is known, then the unknown classes are secure. The classes from which we know at least one element are vulnerable to compromise via Equation 4.

It is sufficient to prove security when exactly two classes are unknown since a procedure for compromising three or more unknown classes could also be used to compromise two unknown classes. Let Y_1 and Y_2 be the two unknown classes, and let y_{11}, \dots, y_{1N_1} be the elements of Y_1 and y_{21}, \dots, y_{2N_2} be the elements of Y_2 . Because data from the known classes is functionally independent of data in the unknown classes, Y_3, \dots, Y_k are irrelevant to compromising Y_1 and Y_2 . Therefore we need only to prove that queries of the form

$$q(z_1 z_2) = a_1 z_1 + a_2 z_2$$

cannot reveal any element of Y_1 or Y_2 .

Suppose that queries $q(z_{i1}, z_{i2})$ have been posed for $i = 1, \dots, m$, yielding responses q_1, \dots, q_m . These queries can be represented by the matrix equation

$$\left[\begin{array}{ccc|ccc} a_{11} & \cdots & a_{1N_1} & b_{11} & \cdots & b_{1N_2} \\ & & & & & \\ & & & & & \\ \hline a_{m1} & \cdots & a_{mN_1} & b_{m1} & \cdots & b_{mN_2} \end{array} \right] \begin{bmatrix} y_{11} \\ \vdots \\ y_{1N_1} \\ \hline y_{21} \\ \vdots \\ y_{2N_2} \end{bmatrix} = \begin{bmatrix} q_1 \\ \vdots \\ q_m \end{bmatrix}$$

where

$$a_{ij} = \begin{cases} a_1, & \text{if } z_{i1} = y_{1j} \\ 0, & \text{otherwise} \end{cases}$$

$$b_{ij} = \begin{cases} a_2, & \text{if } z_{i2} = y_{2j} \\ 0, & \text{otherwise} \end{cases}$$

Now let A_j represent the j^{th} column in the left part of the p by (N_1+N_2) coefficient matrix, and B_j represent the j^{th} column of the right part.

Note that

$$\frac{1}{a_1} \sum_{j=1}^{N_1} A_j - \frac{1}{a_2} \sum_{j=1}^{N_2} B_j = 0.$$

This means that the columns are not independent and that the rank of the coefficient matrix is less than N_1+N_2 -- i.e., the system of equations cannot be solved for any of the N_1+N_2 unknowns.

6. References

- Chi77 Chin, F. Y., "Security in Statistical data bases for queries with small counts," Dept. of Comp. Sci., Univ. of Alberta, 1977.
- Dav76 Davida, B. I. et al., "Data base security," TR-CS-76-14, Dept. of E.E. and Computer Science, Univ. of Wisconsin, July 1976.
- DDL76 DeMillo, R. A., Dobkin, D., and Lipton, R. J., "Even data bases that lie can be compromised," Research Report 67, Yale University, Department Computer Science, May 1976.
- DDS77 Denning, D. E., Denning, P. J., and Schwartz, M. D., "The tracker: a threat to statistical data base security," CSD TR Computer Science Dept. Purdue Univ., Oct. 1977.
- Den78 Denning, D. E., "Are Statistical Data Bases Secure?," Proc. AFIPS NCC, (1978).
- DJL76 Dobkin, D., Jones, A. K., and Lipton, R. J., "Secure data bases: protection against user inference." Research Report 65, Yale University, Department Computer Science, April 1976.
- FeP74 Fellegi, I. P., and Phillips, J. L., "Statistical confidentiality: some theory and applications to data dissemination." Anal. Econ. and Soc'l Measmt. 3, 2 (April 1974), 399-409.
- FoD75 Foster, C. and Daniels, J., "Running by the numbers," Runner's World (July 1975), 14-17.
- Gar68 Garsoux, J. Analyse Mathematique. Dunod, Paris, 1968, p. 444.
- Haq74 Haq, M. I., "Security in a statistical data base." Proc. Am. Soc. for Info. Sci. 11 (1974), 33-39.
- Haq75 Haq, M. I., "Insuring individual's privacy from statistical data base users." Proc. AFIPS (1975), 941-946.
- HoM70 Hoffman, L. J., and Miller, W. F., "Getting a personal dossier from a statistical data bank." Datamation (May 1970), 74-75.
- KaU76 Kam, J. B. and Ullman, J. D., "A Model of Statistical Data Bases and their Security," TR-207, Dept. of EECS, Princeton Univ., June 1976.
- Rei77 Reiss, S. P., "Medians and data base security," Computer Science Dept., Brown Univ., Oct. 1977.
- Sch75 Schlorer, J., "Identification and retrieval of personal records from a statistical data bank," Methods of Information in Medicine 14, 1 (1975), 7-13.

- Sch76 Schlorer, J., "Confidentiality of statistical records: a threat monitoring scheme for on-line dialogue." Methods of Information in Medicine 15, 1 (1976), 36-42.
- Scw77 Schwartz, M. D., "Inference from statistical data bases," Ph.D. Thesis, Computer Science Dept. Purdue University, Aug. 1977.
- SDD77 Schwartz, M. D., Denning, D. E., and Denning, P. J., "Securing Data Bases Under Linear Queries," IFIP Congress Proc., (1977), 395-398.