

On the size of projections:I

E. Gelenbe, D. Gardy

► **To cite this version:**

| E. Gelenbe, D. Gardy. On the size of projections:I. RR-0105, INRIA. 1981. <inria-00076455>

HAL Id: inria-00076455

<https://hal.inria.fr/inria-00076455>

Submitted on 24 May 2006

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

IRIA

CENTRE DE ROCQUENCOURT

Institut National
de Recherche
en Informatique
et en Automatique

Domaine de Voluceau
Rocquencourt
BP 105
78153 Le Chesnay Cedex
France
Tél 954 90 20

Rapports de Recherche

N° 105

ON THE SIZE OF PROJECTIONS : I

Erol GELENBE
Danièle GARDY

Décembre 1981

ON THE SIZE OF PROJECTIONS : I

Erol GELENBE

Danièle GARDY

Abstract

In various applications requiring the processing of large amounts of data, a very common operation is to project a given set of ℓ data points in k -dimensional space into a smaller subspace. The purpose of this note is to estimate the probability distribution of the size of such a projection (i.e. number of data points obtained). Probabilistic assumption concerning the manner in which the initial data is generated are used to obtain the result.

Résumé

Nous nous intéressons au calcul de la taille de la projection d'un nuage de points, donné dans un espace de dimension k , dans un sous-espace. En bases de données il s'agit du problème de la création de vues ou de projections. La solution du problème est donnée dans le cas où tous les nuages initiaux sont équiprobables dans un espace fini, et n'ayant pas de "dépendances fonctionnelles". Un algorithme efficace pour le calcul de la loi de probabilité de la taille des projections est également obtenu.

On the size of projections : I

1. INTRODUCTION

Consider the following problem which arises in various areas of application (physics experiments, census data, etc.) where the collection of large number of data is involved.

The result of the data collection processes is a set

$$T_{\ell k} = \{(t_{11}, \dots, t_{1k}), \dots, (t_{\ell 1}, \dots, t_{\ell k})\}$$

of vectors where each t_{ij} , $1 \leq i \leq \ell$, is an element of the set D_j , $1 \leq j \leq k$. Thus we may consider that T is an ℓ -row and k -column matrix ; the elements of the j -th column all being elements of some set D_j .

For instance, as a result of a physics experiment the mass, position, velocity of one or more particles may be measured. We would then have a table $T_{\ell 3}$ where ℓ is the number of distinct (mass, position, velocity) vectors encountered, the i -th row of $T_{\ell 3}$ being

(mass, position, velocity).

Once such a table has been obtained from the physics experiment it is of interest to compute projections of this table along certain of its columns or coordinates. For instance, one may be interested in obtaining the set of all distinct values of the pair (mass, position). Thus, from $T_{\ell 3}$ we would obtain a new table

$$\pi_3(T_{\ell 3})$$

which would contain ℓ' rows ($\ell' \leq \ell$) and 2 columns, the last (velocity) column being removed.

Just as for T_{ℓ_3} , all of the rows of $\pi_3(T_{\ell_3})$ would be distinct. Thus, for instance, if the measurements yield

$$T_{33} = \begin{bmatrix} 1, 0, 0 \\ 1, 0, 1 \\ 1, 1, 5 \end{bmatrix}$$

we would have

$$\pi_3(T_{33}) = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}$$

The following problem is of interest.

PROBLEM 1 : For a given $T_{\ell k}$ and a given projection π_x along the columns $(1, \dots, x-1, x+1, \dots, k)$ estimate the (size or) number of distinct rows of $\pi_x(T_{\ell k})$.

Let us now consider a seemingly more general problem closely related to this one. The result of the experiment could be a table T_{ℓ_4} , where each row would be

(number, mass, position, velocity)

giving the number of particles counted which have the same (mass, position, velocity) characteristic. We might then be interested in obtaining the total number of particles having the same (mass, position). Thus with the example given above, our initial data might have been

$$T_{34} = \begin{bmatrix} 100, 1, 0, 0 \\ 25, 1, 0, 1 \\ 45, 1, 1, 5 \end{bmatrix}$$

while the result of interest would have been

$$T' = \begin{bmatrix} 125, 1, 0 \\ 45, 1, 1 \end{bmatrix}$$

We see in this case that the number of lines of T' is the same as that of $\pi_3(T_{33})$. In fact the problem of estimating the size of T' is simply an instance of Problem 1, and we see that the estimation of the size of projections is in fact quite general in various problems of data handling.

It is of course also an important issue in data base theory (see for instance [1], [2]). We shall address and solve it in a specific simplified mathematical context.

2. THE FORMAL PROBLEM AND ITS SOLUTION

Let $D_j, 1 \leq j \leq k$, the set of values which may be taken by an element of $T_{\ell k}$, be a finite set.

A table $T_{\ell k}$ is a set of ℓ distinct elements of $D_1 \times D_2 \times \dots \times D_k$. We shall assume that an element $t \in D_1 \times \dots \times D_k$ is generated in the following manner :

$$t = (t_1, \dots, t_k)$$

where t_j is equally likely to be any of the elements of D_j . That is, we treat D_j as a sample space to which we associate a uniform distribution. Furthermore, we assume that t_j is independent of t_m if $j \neq m$.

REMARK 1 : Let $d_i = |D_i|$, and let $h_{\ell k}$ denote the number of distinct tables $T_{\ell k}$. Then

$$h_{\ell k} = \binom{d}{\ell}$$

where $d = d_1 d_2 \dots d_k$. This is simply the number of distinct ℓ -row tables.

RESULT 2 : The number of tables of the form $T_{\ell k}$ whose projection $\pi_j(T_{\ell k})$ along the j -th column contains (exactly) $(\ell - y)$ rows, $0 \leq y \leq \ell - 1$, is

$$(1) \quad Q_{\ell k}^{j,y} = \binom{d/d_j}{\ell-y} \sum_{\substack{n_1, \dots, n_{\ell-y} \geq 0 \\ \sum_{m=1}^{\ell-y} n_m = y}} \prod_{m=1}^{\ell-y} \binom{d_j}{n_m+1}, \quad 1 \leq j \leq k.$$

Proof : Consider the table $T_{\ell-y, k-1}^j$ each of whose $(\ell-y)$ rows have the form

$$(t_1, \dots, t_{j-1}, t_{j+1}, \dots, t_k)$$

There are $h_{\ell-y, k-1}^j$ such tables, where

$$h_{\ell-y, k-1}^j = \binom{d(j)}{\ell-y}, \quad d(j) = d/d_j.$$

To any such table add n_1 replicates of its first row, n_2 of its second row, ..., $n_{\ell-y}$ replicates of its last row, where $n_m \geq 0$ and

$$\sum_{m=1}^{\ell-y} n_m = y$$

to obtain an "intermediate table". Then reconstruct a table of the form $T_{\ell k}$ by introducing the j -th column. For the first $(\ell-y)$ positions of the new j -th column any element of D_j may be used. There are thus $(d_j)^{\ell-y}$ possibilities. Thus for the n_m replicates of the m -th row of the "intermediate table", distinct choices will have to be made : (d_j-1) for the first replicate, (d_j-2) for the second and so on, $d_j - n_m$ for the last replicate.

There are thus

$$\prod_{m=1}^{\ell-y} \binom{d_j}{n_m+1}$$

ways of reconstructing a table of the form $T_{\ell k}$ from a given intermediate table. The total number of tables having k columns and ℓ distinct rows and which yield the same $T_{\ell-y, k-1}^j$ table is therefore

$$\sum_{0 \leq n_1, \dots, n_{\ell-y}} \prod_{m=1}^{\ell-y} \binom{d_j}{n_m+1}$$

$$\sum_{m=1}^{\ell-y} n_m = y$$

hence the result. \square

CONSEQUENCE 3 : The probability that the projection $\pi_j(T_{\ell k})$ will be of dimension (no of lines) x is

$$p_{\ell k}^j(x) = \frac{Q_{\ell k}^{j, \ell-x}}{h_{\ell k}}$$

Proof : It is simply the proportion of tables of the type $T_{\ell k}$ whose projection along the j -th column is of type $T_{x, k-1}^j$. \square

We would now like to evaluate the size of the projection of a table $T_{\ell, k}$ into a subspace composed of the following columns

$$(1, \dots, j_1-1, j_1+1, \dots, j_{s-1}, j_{s+1}, \dots, k)$$

obtained by removing columns (j_1, \dots, j_s) from $T_{\ell k}$. We shall call this projection

$$\pi_{j_1 \dots j_s}(T_{\ell k}) .$$

RESULT 4 : The number of tables of the form $T_{\ell k}$ whose projection $\pi_{j_1 \dots j_s}(T_{\ell k})$ contains exactly $(\ell-y)$ rows, $0 \leq y \leq \ell-1$, is

$$(2) \quad R_{\ell k}^{(j_1, \dots, j_s), y} = \binom{d/(d_{j_1} \dots d_{j_s})}{\ell-y} \sum_{\substack{n_1, \dots, n_{\ell-y} \geq 0 \\ \sum_1^{\ell-y} n_m = y}} \prod_{m=1}^{\ell-y} \binom{d_{j_1} \dots d_{j_s}}{n_m+1}$$

Proof : This is merely a consequence of Result 2. It suffices to notice that the set of tables of the form $T_{\ell k}$ is isomorphic to the set of tables

$$T_{\ell, k-s+1}^{(j_1, \dots, j_s)}$$

obtained by replacing the columns j_1, \dots, j_s of $T_{\ell k}$ by a single column whose elements are chosen from the set $D_{j_1} \times \dots \times D_{j_s}$. The computation of $\pi_{j_1, \dots, j_s}(T_{\ell k})$ is then identical to the projection of

$$T_{\ell, k-s+1}^{(j_1, \dots, j_s)}$$

by the removal of this particular column. □

RESULT 5 : The probability that the projection $\pi_{j_1 \dots j_s}(T_{\ell k})$ is of dimension (no of lines) x is

$$p_{\ell k}^{j_1 \dots j_s(x)} = \frac{R_{\ell, k}^{(j_1, \dots, j_s), \ell-x}}{h_{\ell, k}}$$

3. COMPUTATIONAL ALGORITHMS

Formulae (1) and (2) (the latter having essentially the same form as (1)) are not computationally very efficient. Let us define

$$(3) \quad X_{a,b}^{(v)} \equiv \sum_{\substack{n_1, \dots, n_a \geq 0 \\ \sum_1^a n_m = b}} \prod_{m=1}^a \binom{v}{n_m+1}$$

so that from (1) we have

$$(4) \quad Q_{\ell, k}^{j, y} = \binom{d/d_j}{\ell-y} X_{\ell-y, y}^{(d_j)}$$

and from (2)

$$(5) \quad R_{\ell, k}^{(j_1, \dots, j_x), y} = \binom{d/d_{j_1} \dots d_{j_x}}{\ell-y} X_{\ell-y, y}^{(d_{j_1} d_{j_2} \dots d_{j_x})}$$

Clearly

$$\binom{y+\ell-y-1}{\ell-y-1} = \binom{\ell-1}{\ell-y-1}$$

terms have to be computed and added in order to obtain $X_{\ell,y}^{(\infty,v)}$ using (3), and this can be extremely large even for moderate values of ℓ . For instance, for $\ell = 50$ and $y = 20$ we obtain approximately 2.83×10^{13} which is properly astronomical ! It is therefore useful and even essential to seek a more efficient computational procedure for $X_{\ell,y}^{(v)}$.

Notice that

$$\begin{aligned}
 (6) \quad X_{a,b}^{(v)} &= \sum_{n_a=0}^b \binom{v}{n_a+1} \sum_{\substack{n_1, \dots, n_{a-1} \\ \sum_{i=1}^{a-1} n_i = b - n_a}} \prod_{m=1}^{a-1} \binom{v}{n_m+1} \\
 &= \sum_{z=0}^b \binom{v}{z+1} X_{a-1,b-z}^{(v)} .
 \end{aligned}$$

3.1- Computation of $\{p_{\ell,k}^j(x)\}_{1 \leq x \leq \ell}$

The computation of the probability distribution $\{p_{\ell,k}^j(x)\}_{1 \leq x \leq \ell}$

giving the probability that $\pi_j(T_{\ell k})$ contains x rows requires (see Consequence 3) the computation of $Q_{\ell,k}^{j,\ell-x}$ and hence that of

$$(7) \quad X_{x,\ell-x}^{(d_j)}, \quad 1 \leq x \leq \ell$$

The proposed algorithm is as follows. After setting all $X(d_j) \leftarrow 0$:

$$\begin{aligned}
 &X_{1,0}(d_j) \leftarrow \binom{d_j}{1} ; \quad X_{1,1}(d_j) \leftarrow \binom{d_j}{2} ; \\
 &\underline{\text{for } a=2 \text{ to } \ell \text{ do}} \quad X_{a,0}(d_j) \leftarrow \binom{d_j}{1} * X_{a-1,0}(d_j) \\
 &\underline{\text{for } a=2 \text{ to } \ell \text{ do}} \quad \underline{\text{for } b=1 \text{ to } \ell-a \text{ do}} \quad \underline{\text{for } z=0 \text{ to } b \text{ do}} \\
 &X_{a,b}(d_j) \leftarrow X_{a,b}(d_j) + \binom{d_j}{z+1} * X_{a-1,b-z}(d_j) ; \quad \underline{\text{end all}} .
 \end{aligned}$$

This algorithm will obtain all $X_{a,b}(d_j)$ with $b \leq \ell - a$, and $1 \leq a \leq \ell$.

There are, of course, $\ell^2/2$ values of the pair (x,y) , $1 \leq x \leq \ell$, $0 \leq y \leq \ell-x$, for which (7) has to be evaluated in order to compute the probability distribution $\{p_{\ell,k}^j(x)\}_{1 \leq x \leq \ell}$.

This is obviously because by (1) and (4)

$$(8) \quad p_{\ell,k}^j(x) = \binom{d/d_j}{x} X_{x, \ell-x}^{(d_j)} / h_{\ell,k}$$

However, for any value of (a,b) , $X_{a,b}^{(d_j)}$ is obtained from the values $X_{a-1, b-z}^{(d_j)}$, $0 \leq z \leq b$, in b computational steps. We therefore have a total number of computational steps proportional to

$$\begin{aligned} \sum_{a=1}^{\ell} \sum_{b=0}^{\ell-a} b &= \sum_{a=1}^{\ell} \frac{(\ell-a)(\ell-a+1)}{2} \\ &= \frac{1}{6} \ell(\ell-1)(\ell+1) \end{aligned}$$

4. CONCLUSION

In this note we have examined a problem of interest, and probably of importance in data handling systems or in data base systems : the size of projections of a set of data from a k -dimensional space in which it is given into a smaller subspace. An enumerative approach provides results in the case of uniformly distributed independent samples on a finite dimensional set of possible data values. Certain extensions, in particular to dynamically varying data sets and to certain cases of dependence (e.g. "functional dependence" in data bases), will be treated in subsequent papers.

ACKNOWLEDGEMENTS

This work was motivated by discussions with W.A.ARMSTRONG and with Ph. RICHARD.

- [1] Ph. RICHARD - Thèse de Doctorat de 3ème cycle, Orsay, 1980.
- [2] Ph. RICHARD - "Evaluation of the size of a query expressed in relational algebra", Proceedings ACM-SIGMOD International Conference on Management of Data, April 1981, pp. 155-163.

