

Creation of Information Profiles in Distributed Databases as a n -Person Game

Juliusz L. Kulikowski

Institute of Biocybernetics and Biomedical Engineering PAS, Poland.

jlkulik@ibib.waw.pl

Abstract. It is considered a problem of forming information profiles of distributed data- and/or knowledge bases as a result of satisfying the requirements representing the information profiles of the customers. It is shown that the interests of data- and knowledge bases managers are not fully convergent and that they participate in a composite, partially co-operative, partially non-co-operative game. There is given a formal description of the strategies used in this game, as well as the methods of decision making of the players.

1 Introduction

A spontaneous and dramatic development of open-access data- and knowledge bases in the last decades has led to a new situation in many areas of human activity depending on the access to information resources. This remark concerns various branches of scientific research, technology, education, administration, trade, health services, national security, natural environment protection, etc. For the existence and development of all the above mentioned areas access to information resources satisfying specific requirements of the users is necessary. Distributed databases accessible through Internet (or through any other computer networks) make possible reaching higher quality of human activities and, as a result, they open a new era in modern civilisation development. But at the same time their impact on our life has caused new problems that existed never before or existed only in a germinal state. Till databases were created by the same organisations they were dedicated to the interests of database managers were close to those of database users and the area of conflicts between them was strongly limited. New situation arose when government or other higher administrative authorities tried to initialise design and construction of computer-based information systems in order to force higher effectiveness in sub-ordered organisations. The goals of information systems' sponsors, designers and users were divergent and, as a result, many so-designed information systems failed as being not accepted by their potential users. Free information market gives another possibility of information systems creation according to the requirements of various subjects. In this case existence of no general "optimum" is assumed; information creators, systems managers, and information users can express their proper goals and they participate in a multi-person game trying to reach their individual optima. Simultaneous reaching of all so-defined optima is impossible; however, the market mechanisms make possible reaching a common balance point being a compromise between the partners' expectations.

General concepts of the theory of games we owe to J. von Neumann and O. Morgenstern [1]. For more than fifty years various types of games: two- and multi-person, discrete- and continuous-strategy, differential, antagonistic and non-antagonistic, co-operative and non-co-operative, one- and multi-level, etc. were investigated [2,3,7]. Various areas of game theory applications: in business, technology, military service, etc., were also considered. The work [4] by E.A. Berzin where game theory application to distribution of resources has been investigated is close to our interests. Some aspects of games played in information systems design and maintaining were also mentioned in [5] while in [6] the role of self-organisation in distributed databases development was considered. The aim of this paper is presentation of a more detailed model of distributed information resources gathering and their profiles forming based on the theory of n -person games.

2 Basic Model Assumptions

It will be considered a system of primary information supply (IS) and a one of information distribution (ID) as two complementary components of information market. The IS system is the one where information in the form of various types of documents (electronic, photo, multi-media, hard copy manuscripts, publications, etc.) is offered to the customers. For many years this type of information distribution and exchange was prevailing.

This situation has been changed with computer networks development; the action of ID system is based on electronic documents exchange through computer networks. ID system thus consists of a set of open-access data banks ($OADB$) and of a number of dedicated local data banks (LDB), mutually connected by computer network, as shown in Fig. 1.

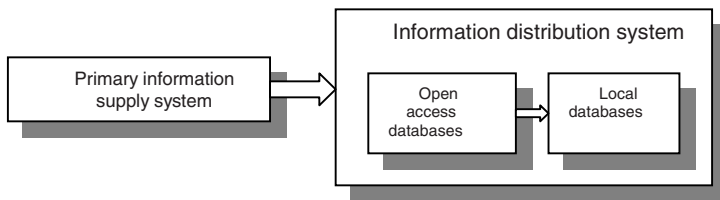


Fig. 1. A model of information storage and distribution

Information resources of $OADB$ s are permanently extended and supplied with documents from the IS system. The role of $OADB$ s consists in brokerage: they buy information documents, select information from them, transform, store and sell or distribute it among the LDB s, the last being explored by organisations or by individual users. Information resources of the LDB s are thus completed according to the needs of the users. The LDB s can be supplied by information not only from $OADB$ s, but also from their proper information sources, as well as directly from the IS system.

It is assumed that each document offered at the information market can be described by the following characteristics: 1) *formal properties*: type of the

document, its author(s), title, volume, editor, date of edition, etc.; 2) *contents* (described by keywords, descriptors, etc.); 3) *acquisition conditions* (name of the seller, price, forms of payment, delivery time, number of available copies of the document, etc.); 4) *supply indicator* (number of available copies of the document, if limited, infinite otherwise). The document characteristics can be formally represented as elements of a *document characteristics space DCS* being a Cartesian product:

$$DCS = C_f \times C_c \times C_a \times C_s \tag{1}$$

where C_f , C_c , C_a and C_s stand, correspondingly, for the sets of formal, contents, acquisition and supply characteristics' possible values. We shall denote by \mathbf{h}_i , $\mathbf{h}_i \in DCS$, the characteristic of a document \mathbf{x}_i , $i \in [1,2,3,\dots]$.

The state of information market is a time-varying process: at any fixed time-instant t_0 , $t_0 \in T$ (T being a real discrete time-axis), the instant-value of the process is given by a finite subset $X(t_0) \subset DCS$ of the documents that actually are offered on the market. The members \mathbf{x}_i of $X(t)$ thus arise at certain time-instants and next, after a certain time-delay they may disappear. However, formal description of a time-process' values in the set theory terms is not very convenient. Instead of this it will be taken into account a set of time-functions: ξ_i :

$$T \rightarrow DCS \tag{2}$$

assigning to each time-instant $t \in T$ a document characteristic (an element of DCS) so that the components of C_f and C_c remain constant while those of C_a and C_s may be varying in time. The subsets $X(t)$ contain only the elements (documents) for which the indices δ_i of the components of C_a take positive, non-zero values. For those elements the components of C_a are also defined, otherwise the values \emptyset (= *undefined*) to them are assigned. However, the vector $\xi(t)$ consisting of linearly ordered components $\xi_i(t)$ is not known but for the past and present time-instants t only. It is reasonable to consider it as a realisation of a stochastic vector process $\Xi(t)$ describing the states of the information market changing in time. The subsets $X(t)$ determine, at the same time, the areas of possible *instant decisions* of the *OADB* managers. Their aim consists in actualisation and extension of the *OADB*s according to the expected demands of the customers. Therefore, at any t they can chose the following decisions: 1^o to select and to acquire new documents in order to include and to keep them in the data banks, 2^o to select some documents and to delay a final decision about their acquisition, and 3^o to reject all actual proposals concerning documents selling. In the cases 1^o and 2^o the decision made by the ν -th *OADB* manager ($\nu = 1,2,\dots,N$, where N denotes the total number of *OADB* managers) takes the form of a subset:

$$\psi_\nu(t) \subseteq DCS \tag{3}$$

whose members correspond to the selected elements of $X(t)$ and are such that:

1^o the projections of the members of $\psi_\nu(t)$ on $C_f \times C_c \times C_a$ are equal to the corresponding elements of $X(t)$ projected on $C_f \times C_c \times C_a$;

2^o the values of the components β_i of C_s in the members of $\psi_\nu(t)$ satisfy the inequalities

$$0 \leq \beta_i(t) \leq \delta_i(t), \tag{4}$$

which means that the documents acquisition requirement can not exceed the corresponding supply indicators. The document suppliers then collect the requirements and try to realise them. The number of sold document's copies can not exceed the declared supply indicator. Therefore, it may happen that some document acquisition requirements are not satisfied. However, the strategy of clients' selection on the information market here will not be considered. In any case, when coming to the next $(t+1)$ time-instant the supply indicators should be actualised: reduced by the numbers of sold document copies and increased by the numbers of the new-supplied ones.

The documents acquired by the *OADB*s can be included into the data banks directly or after a transformation process changing their form or generating some secondary documents on the basis of information selected from the primary ones. In similar way as shown before a modified *document characteristics space* can be defined:

$$DCS^{(\nu)} = C^{(\nu)}_f \times C^{(\nu)}_c \times C^{(\nu)}_a \quad (5)$$

Here the notions $C^{(\nu)}_f$, $C^{(\nu)}_c$ and $C^{(\nu)}_a$ are, in general, similar to the C_f , C_c and C_a used in (1); however, a dependence on ν indicates that each $(\nu$ -th) *OADB* may use its proper language and standards for data files characterisation and define data acquirement conditions for the customers. So, a subset $K^{(\nu)} \subseteq DCS^{(\nu)}$ plays the role of a catalogue of data offered to the managers of *LDB*s or to individual users. A projection of $K^{(\nu)}$ on the document characteristic subspace $C^{(\nu)}_f \times C^{(\nu)}_c$, $L^{(\nu)} \subseteq C^{(\nu)}_f \times C^{(\nu)}_c$, will be called a *profile of the ν -th OADB* ($OADB^{(\nu)}$), and its elements will be denoted by $\lambda^{(\nu)}$. A substantial difference between the formerly defined characteristics h_i and the $\lambda^{(\nu)}$ ones consists in the fact that h_i describes the document in its original form while $\lambda^{(\nu)}$ describes secondary documents in electronic form. Forming the profiles $L^{(\nu)}$ of the *OADB*s is the main element of long-term strategies of *OADB*'s managers.

Then, let us take into consideration the managers' of *LDB*s point of view. They represent the interests of some groups of information users (or are information users themselves). The users need to have easy access to information resources suitable for satisfying their intellectual (educational, cultural, etc.) interests or for solving some professional (technological, administrative, etc.) tasks. Let us assign index μ to a certain group of information users. Then their information needs can be formally represented by subsets of a Cartesian product describing an *information requirements space*:

$$IRS^{(\mu)} = C^{*(\mu)}_f \times C^{*(\mu)}_c \times C^{*(\mu)}_a \quad (6)$$

Once more, the notions $C^{*(\mu)}_f$, $C^{*(\mu)}_c$ and $C^{*(\mu)}_a$ are, in general, similar to the C_f , C_c and C_a ones; however, a dependence on m shows that each $(m$ -th) *LDB* may use its proper language and standards for data files characterisation and define additional conditions for data acquisition (like admissible cost, time-delay, etc.). In particular, the sets $C^{*(\mu)}_f$, $C^{*(\mu)}_c$ and $C^{*(\mu)}_a$ may contain an element $*$ (*any possible*) to be used if some characteristics of data files or records are not fixed by the users. The subset $R^{(\mu)} \subseteq IRS^{(\mu)}$ characterising information needs of the user(s) can be used in two ways: 1) for information retrieval in $LDB^{(\mu)}$, and 2) for actualisation of *local information resources* ($LIR^{(\mu)}$) in $LDB^{(\mu)}$.

$LIR^{(u)}$ is also a subset of $IRS^{(u)}$. An information retrieval order can be realised if a certain *consistency level* between $R^{(u)}$ and $LIR^{(u)}$ is reached, otherwise it is necessary to import the necessary data from the *OADB*'s. However, the managers of *LDB*s may conduct a more active policy of users' requirements realisation. This can be reached by a permanent monitoring of the information requirements flow: ... $R^{(u)}(t-3)$, $R^{(u)}(t-2)$, $R^{(u)}(t-1)$, $R^{(u)}(t)$ (t being the current time) in order to define a preferable *local resources profile*. The last for the given $LDB^{(u)}$ can be defined as a subset

$$\Lambda^{(u)} \subseteq C_f^{*(u)} \times C_c^{*(u)} \tag{7}$$

such that if $LIR^{(u)} = \Lambda^{(u)}$ then a considerable part of expected information requirements can be directly realised. The managers of *LDB*s try to achieve this situation within their possibilities as it will be shown below.

3 Strategies for *OADB* Managers

The *OADB*s' and *LDB*s' managers are interested in realisation of information requirements of their customers. The managers of *OADB*s are customers on the *IS* market and, at the same time, they are data sellers with respect to the managers of *LDB*s. On the other hand, the last ones are data suppliers for the information users. Buying a document x_i available on *IS* market needs covering a cost κ_i being indicated as a component of the corresponding document characteristic. The same document (or data drawn from it) included into the information resources of $OADB^{(v)}$ is expected to be distributed among a number of *LDB*s to the information profiles of which it suits. As a consequence, the manager of $OADB^{(v)}$ expects to reach proceeds of $r_i^{(v)}$. Finally, its expected profit from buying and distributing x_i is

$$c_i^{(v)} = r_i^{(v)} - \kappa_i \tag{8}$$

It might seem that there is a simple decision rule for the $OADB^{(v)}$ manager:

$$\left. \begin{array}{l} \text{buy } x_i \text{ if } c_i^{(v)} \geq e_i^{(v)} > 0, \\ \text{delay the decision if } 0 < c_i^{(v)} < e_i^{(v)} \\ \text{do not buy it otherwise.} \end{array} \right\} \tag{9}$$

where $e_i^{(v)}$ is a threshold. However, there arise the following problems: 1° how to evaluate $c_i^{(v)}$ (or $r_i^{(v)}$, see (8)), and 2° how to establish $e_i^{(v)}$?

For answering the first question there can be made the following assumptions:

a/ expected proceeds $r_i^{(v)}$ can be described by an increasing function of the number of customers whose information profiles $\Lambda^{(u)}$ are *consistent* with the information characteristics h_i of x_i and of a total measure of this *consistency*;

b/ the information profiles $\Lambda^{(u)}$ of the *LDB*s are not known exactly to the manager of $OADB^{(v)}$, he can only approximate them in the form of information profile $L^{(v)}$ constructed on the basis of all past and actual information requirements from the *LDB*s;

c/ $r_i^{(v)}$ is a decreasing function of other *OADB*s in the *ID* system that will offer x_i and/or some secondary information drawn from it.

Then, it arises the next problem: how to evaluate the *measure of consistency* between a document characteristic \mathbf{h}_i and the profile of $\mathbf{L}^{(v)}$, assuming that (3) holds and $\mathbf{h}_i \in C_f^{(v)} \times C_c^{(v)}$. Let us remark that in general the elements of $C_f^{(v)} \times C_c^{(v)}$ are not vectors in algebraic sense, but rather some strings of elementary data of various formal nature. Therefore, it is not possible to take an ordinary *distance measure* concept as a basis of a *consistency measure* definition. However, the last can be based on a generalised, *multi-aspect similarity measure* concept proposed in [8]. In this case, if A is a non-empty set, then a similarity measure between its elements can be defined as a function:

$$\sigma: A \times A \rightarrow [0,1]_c \quad (10)$$

where $[p,q]_c$ denotes a continuous interval between p and q , such that: $a/$ for each $a \in A$ there is $\sigma(a,a) \equiv 1$; $b/$ for any $a, b \in A$ there is $\sigma(a,b) \equiv \sigma(b,a)$; $c/$ for any $a, b, c \in A$ there is $\sigma(a,c) \geq \sigma(a,b) \cdot \sigma(b,c)$.

If $\mathbf{f}^{(r)} = [f_1^{(r)}, f_2^{(r)}, \dots, f_p^{(r)}, \dots, f_g^{(r)}]$ and $\mathbf{f}^{(s)} = [f_1^{(s)}, f_2^{(s)}, \dots, f_p^{(s)}, \dots, f_g^{(s)}]$ are two strings of characteristics whose components are of various formal nature then a measure of multi-aspect similarity can be defined as a product:

$$\sigma(\mathbf{f}^{(r)}, \mathbf{f}^{(s)}) = \sigma_1(f_1^{(r)}, f_1^{(s)}) \cdot \sigma_2(f_2^{(r)}, f_2^{(s)}) \cdot \dots \cdot \sigma_g(f_g^{(r)}, f_g^{(s)}) \quad (11)$$

This definition can be used directly to the similarity evaluation of documents characteristics. Let $\sigma(\mathbf{h}_i, \mathbf{h}_j)$ be a similarity measure described on the Cartesian product $\mathbf{A} = C_f^{(v)} \times C_c^{(v)}$. Then it will be said that a member $\mathbf{h}_p, \mathbf{h}_j \in \mathbf{A}$ is *adherent* to a subset $\mathbf{L}^{(v)} \subseteq \mathbf{A}$ on a level ε , $0 < \varepsilon \leq 1$, if there is at least one element $\mathbf{h}_j \in \mathbf{A}$ such that $\varepsilon \leq \sigma(\mathbf{h}_p, \mathbf{h}_j) \leq 1$. Here ε is a threshold chosen according to the application requirements.

Adherence of \mathbf{h}_i to $\mathbf{L}^{(v)}$ on a fixed level is necessary, but not a sufficient condition for making a positive decision according to the rule (9). For this purpose the *OADB*^(v) manager should also take into account: 1° how many characteristics \mathbf{h}_j in $\mathbf{L}^{(v)}$, for all possible \mathbf{h}_p , satisfy the adherence condition, 2° how many customers have declared their interests in acquiring the data characterised by \mathbf{h}_p , and 3° how long time has been passed since the last call for \mathbf{h}_p . The corresponding, additional data can be stored and included into the information requirement characteristics (see (6)). Taking them into account a *measure of consistency* $\Gamma(\mathbf{h}_i, \mathbf{L}^{(v)})$ between \mathbf{h}_i and $\mathbf{L}^{(v)}$ can be defined as it will be illustrated below.

Let us assume that the contents of documents are characterised by keywords that are presented in a linear order: w_1, w_2, w_3, \dots , etc. Let us take into account two documents, whose characteristics \mathbf{h}_p and \mathbf{h}_j contain, correspondingly, the subsets of keywords W_i and W_j . There will be considered the sets: $W_i \cup W_j$ and $W_i \cap W_j$. If the cardinal number of a set W is denoted by $|W|$ then the similarity measure of the above-mentioned sets can be defined as

$$\sigma(W_i, W_j) = \frac{|W_i \cap W_j|}{|W_i \cup W_j|} \quad (12)$$

This similarity measure can be taken as a basis of the $\Gamma(\mathbf{h}_i, \mathbf{L}^{(v)})$ definition. First, let us remark that $\mathbf{L}^{(v)}$ can be formally interpreted as a virtual document consisting of all documents whose characteristics were in the past required by a considerable part of customers. Therefore, the formula (12) can be used for evaluation of similarity

between h_i and $L^{(v)}$, as well. Let us suppose that in the given $OADB^{(v)}$ the requirements are registered in such a way that to any keyword w_α there are assigned the numbers $e_{\alpha,\tau}$, $\tau = 0, 1, 2, \dots, T$ indicating how many times w_α was mentioned in the information calls in the present ($\tau = 0$), as well as in the former time-periods (years). Then the following weight coefficient:

$$\lambda_\alpha = \sum_{\tau=0}^T \frac{e_{\alpha,\tau}}{\tau + 1} \tag{13}$$

and, at last we can put

$$\Gamma(h_i, L^{(v)}) = \sigma(h_i, L^{(v)}) \cdot \lambda_\alpha \tag{14}$$

where the sum λ_α on the right side is taken over all keywords occurring both in h_i and $L^{(v)}$. The manager of $OADB^{(v)}$ thus can assume a proportionality:

$$c_i^{(v)} = k \cdot \Gamma(h_i, L^{(v)}) \tag{15}$$

(k being a positive coefficient of proportionality) meaning that the greater is the consistency measure between h_i and the profile $L^{(v)}$, the higher are the expected proceeds from selling information drawn from x_i . However, in the situation when the supply of x_i on the market is high (i.e. > 1) the same document can be bought and distributed by other $OADB$ s. In such case, assuming that the number of LDB s interested in acquiring x_i is unchanged, the incomes of the $OADB$ s will be inversely proportional to the number of $OADB$ s acquiring x_i . Therefore, the managers of $OADB$ s in this case play a competitive game with fixed positive total gain. They may all acquire the given document; in such case, if they do it at the same time, they will share the gain proportionally to the number of customers that will require access to data from the document. However, the number of $OADB$ s buying the document a priori is not known. In such case the manager of $OADB^{(v)}$ may, first, delay his decision in order to recognise the actions of other players. However, if the delay is large then the expected proceeds will be reduced as well.

The risk of acquiring documents which will not be needed by the customers may be not accessible for the $OADB$ s' managers. In such case they may change the game rules by establishing a coalition among them. The coalition consists in specialization of their information resources so that the information profiles $L^{(v)}$ are (at least, approximately) disjoint. As a consequence, the $OADB$ s will share the information requirements of the customers according to their contents. In such case the expected incomes of the $OADB$ s lose the factor of uncertainty connected with the number of competitive $OADB$ s offering access to the same documents.

4 Strategies for LDB Managers

The relationships between the LDB managers and the information users is not the same as this one between the $OADB$ s' and LDB s' managers. The difference consists in the fact that the LDB s' managers and their information users usually act within the same organisations. Therefore, up to a certain degree their interests are convergent. However, a source of conflicts may be connected with the fact that the LDB s'

manager has at his disposal a limited amount of financial resources that can be used to a realisation of information requirements of his clients. In such case he must establish a system of priorities in acquiring new data from the *OADB*s. As a consequence, the users attached to a certain *LDB*^(μ) may have partially common and partially competitive interests. Assuming that the *LDB*s' manager has no reason for distinguishing some information users against some other ones he will try to make his decisions on basis of the information profile $A^{(\mu)}$. His problem is thus as follows: 1) He has at his disposal an amount ζ of financial means that he can spend for extending the *LIR*^(μ) for a certain time T^* ; 2) He knows the actual information requirements $R^{(\mu)}$ of his clients (data users) that can not be realised on the basis of *LIR*^(μ); 3) He knows the actual information resources of the *OADB*s as well as the corresponding data selling and/or data access conditions.

How to realise the information needs of the clients within the financial and/or formal limits?

It is assumed that the prices of an incidental access to some data and of data buying for their permanent using within the organisation are different (the last being usually much higher). Therefore, if an analysis of the *OADB*s' resources shows that a certain requirement can be realised by importing data then it is necessary to determine the expected costs of: a / data buying and b / data access during the time-interval T^* for various *OADB*s and to find out the minimum ones. At the next step the conflict of interests of various clients should be solved: it has been determined a set of data acquisition offers: O_1, O_2, O_3, \dots etc., each one being characterised by its expected cost χ_i and by the subset of clients that might be interested in using the given data in the time-interval T^* . The problem is: how to select a subset of the offers for their final acceptance?

The problem can be solved if a system of ordering the offers is defined. Each offer can be characterised by a vector whose components indicate: 1° the expected cost χ_i , 2° the number b_i of clients that are interested in using the given data, and 3° the mean measure of similarity s_i of the offered data to the actual (and expected, if possible) data requirements. Therefore, it arises the problem of semi-ordering of the vectors $u_i = [\chi_i, b_i, s_i]$ in a three-dimensional space. The problem can be solved easily on the basis of K -space (Kantorovich space) concept [9]. Then the strategy of the *LDB* manager consists in rearranging the offers: $O_{\rho_1}, O_{\rho_2}, O_{\rho_3}, \dots$ etc. ($\rho_1, \rho_2, \rho_3, \dots$ etc. being some integers) and in accepting for realisation as much offers as possible within the financial limits.

5 Conclusions

We tried to show that the relationships between the managers of *OADB*s and of *LDB*s have the form of a n -person game of partially co-operative, partially non-co-operative type with incomplete information. It is possible to define the strategies of the players. The strategies are realised by creation of database profiles that are used in making decisions concerning documents or data acquisition. The decision rules of the players can be strongly optimised only in particular, rather simple cases. In general, it has been shown that a co-operation between the database managers may them lead to less risk in reaching their goals.

References

1. Neumann, von J., Morgenstern O.: Theory of Games and Economic Behavior. Princeton (1944)
2. Dubin, G.N., Suzdal, V.G.: Introduction to Applied Theory of Games. Nauka, Moskva (1981) (in Russian)
3. Owen, G.: Game Theory. W.B. Saunders Company, Philadelphia (1968)
4. Berzin, E.A.: Optimum Distribution of Resources and Theory of Games. Radio i Svjaz', Moskva (1983) (in Russian).
5. Kulikowski, J.L.: Problems of Information Systems Design in Social Environment (in Polish). Techniki i metody rozproszonego przetwarzania danych, cz. V (M. Bazewicz ,ed.). Wroclaw Univ. of Tech. Press (1990).
6. Kulikowski, J.L.: Self-organization in Distributed Data Bases. Information Systems' Architecture and Technology ISAT'93 (M. Bazewicz ed.). Wroclaw Univ. of Tech. Press (1993), 94-104
7. Vorob'ev, N.N.: Principles of game theory. Non-co-operative games. Nauka, Moskva (1994)
8. Kulikowski, J.L.: From Pattern Recognition to Image Interpretation. Biocybernetics and Biomedical Engineering, vol. 22, No 2-3 (2002) 177-197.
9. Kantorovich, L.V., Vulich B.Z., Pinsker A.G.: Funkcjonalnyj analiz v poluuporiadocennyh prostranstvach. GITTL, Moskva (1959).