b 1426/101
1890 1

Zdzisław Pawlak

# Mathematical foundations
# of information retrieval

Zdzisław Pawlak

# MATHEMATICAL FOUNDATIONS OF INFORMATION RETRIEVAL
## (A new approach)

# 101

Mailing address:   Prof. Dr. Z. Pawlak
                   Warszawa
                   ul. Gamerskiego 3 m.5
                   POLAND

A b s t r a c t

This note contains a simple mathematical formulation of
basic ideas concerning information retrieval and its computer
implementation. The presented theory is based on the results
given in [1], [2] and [3].


С о д е р ж а н и е

Работа касается математической формулировки основных тер-
минов связанных с проблемой поиска информации и применения
этого аппарата в цифровых машинах. Представленная теория опи-
рается на результаты даны в [1], [2] и [3].


S t r e s z c z e n i e

Praca dotyczy matematycznego sformułowania podstawowych
pojęć związanych z problemem wyszukiwania informacji oraz zasto-
sowania tego aparatu na maszynach cyfrowych. Przedstawiona
teoria oparta jest na wynikach podanych w [1], [2] i [3].

This note contains a simple mathematical formulation of basic ideas concerning information retrieval and its computer implementation. The presented theory is based on the results given in [1], [2] and [3].

### 1. Descriptive systems

By a descriptive system we mean triplet $D = \langle A_D, X_D, \sigma_D \rangle$ (or briefly $D = \langle A, X, \sigma \rangle$ ), where

A  - is a (finite or infinite) set; elements of A are
    called objects of $D$ ,

X  - is a finite set of symbols; elements of X are
    referred to as elementary descriptors of $D$ ,

$\sigma \subseteq A \times X$ - is a binary relation, called description
    relation (or description) in $D$ .

Relation $\sigma$ may be replaced by the function:

$$\psi : X \to 2^A$$

such that:

$$\psi(x) = \left\{ a \in A ; \; \sigma(a,x) \right\} .$$

Let $X^*$ denote the smallest set containing X and such that: if $x, y \in X^*$ then $x \wedge y$, $x \vee y \sim x$ are also in $X$ , and let $\psi^*$ be defined as follows:

$$\Psi^*(x) = \Psi(x) \text{ if } x \in X$$

$$\bigwedge_{x,y \in X^*} \Psi^*_{(x \vee y)} = \Psi^*(x) \cup \Psi^*(y)$$

$$\bigwedge_{x,y \in X^*} \Psi^*_{(x \wedge y)} = \Psi^*(x) \cap \Psi^*(y)$$

$$\bigwedge_{x \in X^*} \Psi^*(\sim x) = A - \Psi^*(x),$$

where $\Psi^*$ is the extension of $\Psi$ (in what follows, the asterisk (*) will be omitted.)

Descriptors $x, y \in X^*$ are said to be equal ($x \equiv y$) iff

$$\Psi(x) = \Psi(y)$$

If $x \equiv y$ does not hold, then we say that $x$ is different from $y$. We assume that all elementary descriptors in $\mathcal{D}$ are always different.

Theorem 1. For every descriptive system $\mathcal{D} = \langle A, X, \delta \rangle$, the number of different descriptors is finite and is not greater than $2^{2^{\overline{\overline{X}}}}$.

## 2. Atomic descriptors

Every product of all elementary descriptors in $\mathcal{D}$ with or without negation

$$x_1^{i_1} \wedge x_2^{i_2} \wedge \ldots \wedge x_k^{i_k}, \qquad x_j^{i_j} \in X$$

where $i_j = 0$ or $1$, $k = \overline{\overline{X}}$, and

$$x_j^0 = x_j, \quad x_j^1 = \sim x_j$$

will be called atomic descriptor in $\mathcal{D}$.

Of course for every $\mathcal{D}$ there are at most $2^{\overline{\overline{X}}}$ different atomic descriptors in $\mathcal{D}$.

If $x$ is an atomic descriptor in $\mathcal{D}$ then $\Psi(x)$ is called atom in $\mathcal{D}$.

Descriptors $x, y \in X^*$ are said to be independent iff

$$\Psi(x) \cap \Psi(y) = \overline{\phi}.$$

Theorem 2. Every two different atomic descriptors in $\mathcal{D}$ are independent.

Theorem 3.

$$\bigcup_{x \in \overline{X}_{\mathcal{D}}} \Psi(x) = A_{\mathcal{D}}$$

where $\overline{X}_{\mathcal{D}}$ denotes the set of all atomic descriptors in $\mathcal{D}$.

Theorem 4. Every elementary descriptor $x \in X_{\mathcal{D}}$ may be represented as

$$x = x_1 \vee x_2 \vee \ldots \vee x_n, \quad x_i \in \overline{X}_{\mathcal{D}},$$

where $x_1, x_2, \ldots, x_n$ are all atomic descriptors in $\mathcal{D}$ containing $x$.

Theorem 5. Every descriptor $x \in X_{\mathcal{D}}^*$ may be represented as the sum of some atomic descriptors in $\mathcal{D}$.

By means of the theorem 5 we are able to represent descriptors in some standard (normal) form.

## 3. Remarks on implementation

Given some set of objects A (for example books, papers, documents, etc.) and a set of elementary descriptors X (for example, authors' names, languages, key words, etc.) thus the relation $\delta$ is defined. Now we may ask about some sets of

objects defined by any compound descriptor $x \in X_{\mathcal{D}}^{*}$.

The set $B \in A_{\mathcal{D}}$ will be called <u>descriptive</u> in $\mathcal{D}$ iff there exists $x \in X_{\mathcal{D}}^{*}$ such that

$$\Psi_{(x)} = B$$

Let $\mathcal{D}(A)$ denote the set of all descriptive sets in $\mathcal{D}$ .

From theorem 1 follows that there is only a finite number of descriptive sets in any $\mathcal{D}$ . So we are unable to "describe" (in the sense of this paper) all subsets of $2^{A}$ in $\mathcal{D}$. By theorem 5 it follows that only the sets which are sums of atoms are descriptive in $\mathcal{D}$ .

This result leads to very simple computer implementations of any information retrieval system:

given any descriptor $x \in X_{\mathcal{D}}^{*}$ by means of theorem 4 we can represent it in a normal form, and then find out the corresponding atoms and form its sum.

This method of searching for sets of documents satisfying some conditions leads to very quick, efficient and simple computer algorithms. Moreover, in the case when the set of elementary descriptors should be extended, the whole system may be easily extended too without destroying and rebuilding already working system.

Many thanks due to Dr. A. Mazurkiewicz for valuable comments.

## Literature

[1] A. Mostowski, K. Kuratowski: Teoria mnogości, PWN, 1966

[2] Z. Pawlak: About the meaning of personal pronouns (to appear in Cahiers de Linguistique Théoretique et Appliquée, Vol. X, 1973, No 1)

[3] Z. Semadeni: Logical kits, Manuscript, 1971.