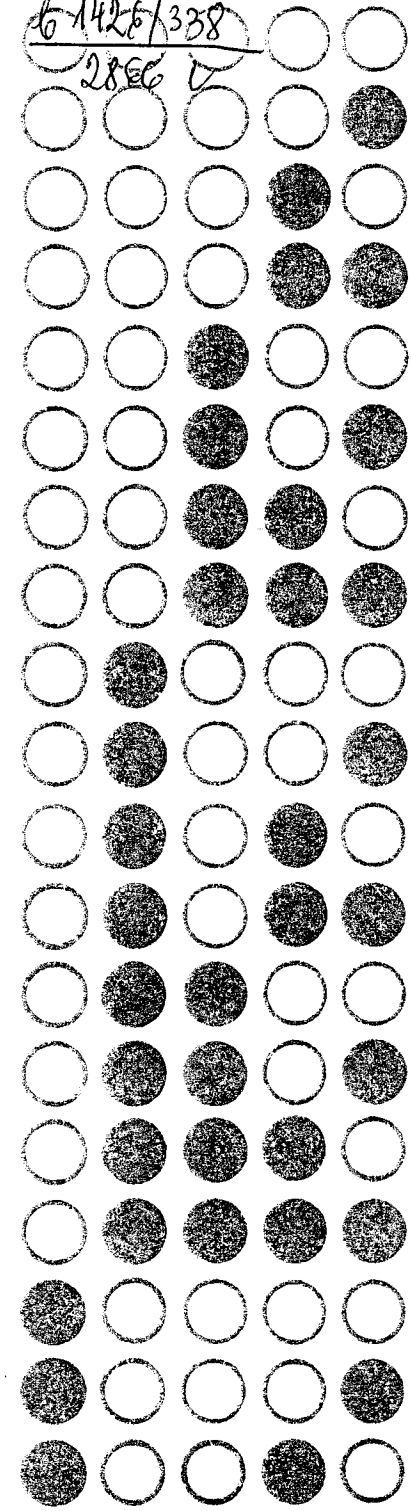


6-1426/338

2886



Zdzisław Pawlak

Information systems

338

1978

WARSZAWA

Zdzisław Pawlak

INFORMATION SYSTEMS

338

Warsaw 1978

R a d a R e d a k c y j n a

A. Blikle (przewodniczący), S. Bylka, J. Lipski (sekretarz),
L. Łukaszewicz, R. Marczyński, A. Mazurkiewicz, T. Nowicki,
Z. Pawlak, D. Sikorski, Z. Szoda, M. Warmus (zastępca prze-
wodniczącego)



Mailing address: Prof. dr Zdzisław Pawlak
Institute of Computer Science
Polish Academy of Sciences
P.O. Box 22
00-901 Warsaw
Poland

sygn.

b 1426/338

nr inw.

2866 L

Printed as a manuscript

На правах рукописи

Nakład 700 egz. Ark. wyd. 0,5; ark. druk. 1,00.
Papier offset. kl. III, 70 g, 70 x 100. Oddano do
druku we wrześniu 1978 r. W. D. N. Zam. nr 618/78

S-83

Abstract . Содержание . Streszczenie

This is new simplified version of ideas given in: W. Marek,
Z. Pawlak: Information Storage and Retrieval Systems. Mathema-
tical Foundations, Theoretical Computer Science 2(1978) pp.
331-354, and in Z. Pawlak: Mathematical Foundations of Infor-
mation Retrieval, CC PAS Reports No 101(1973), concerning infor-
mation systems.

Presented approach offers more adequate tools for dealing
with problems related to information systems. In particular,
the notion of independence of attributes is defined, as well as
minimization problem for attributes is formulated.

Информационные системы

Приводится новая версия формализации понятия информаци-
онной системы, рассматриваемой в работах В.Марка и З.Павлика
"Information Storage and Retrieval Systems. Mathematical
Foundations". Theoretical Computer Sciences 2 /1978/, pp 331 -
354, и З.Павлика "Mathematical Foundation of Information
Retrieval". CC PAS Reports No 101 /1973/, которая дает возмож-
ность введения простой формулировки проблем, связанных с дан-
ной областью. В частности уточняется понятие зависимости при-
знаков и минимизации множества признаков.

Systemy informacyjne

Podano nową wersję formalizacji pojęcia systemu informacyjnego rozpatrywanego w pracach W. Marka i Z. Pawlaka "Information Storage and Retrieval Systems. Mathematical Foundations". Theoretical Computer Sciences 2 (1978) pp. 331-354 i Z. Pawlaka "Mathematical Foundation of Information Retrieval". CC PAS Reports No 101 (1973), która pozwala na prostsze niż poprzednio formułowanie problemów z tego zakresu. W szczególności sprecyzowano pojęcie zależności atrybutów oraz minimalizacji zbioru atrybutów.

This is somewhat modified version of ideas given in [1] and [2].

1. INFORMATION SYSTEM

Definition. An information system is a 4-tuple

$$S = \langle X, A, Q, \xi \rangle,$$

where:

X - is a set of objects,

A - is a set of attributes,

Q - is a set of descriptors,

ξ - is an information function from $X \times A$ into

Let $Q_a \subset Q$, $a \in A$, denote subset of Q defined as follows:

$$Q_a = \{q \in Q : \forall x \xi(x, a) = q\}$$

Thus Q_a is the set of all descriptors assigned to the attribute a. We shall also refer to the elements of Q_a as values of a.

The following notation will be used in the paper:

- a) $\text{Inf}(S) = Q^A$ - the set of all informations in S
- b) For x in X, ξ_x is to mean a function from A into Q such that $Q_x(a) = q$ if and only if $\xi(x, a) = q$

We refer to \mathcal{I}_x as the information about x in S

c) For $\varphi \in \text{Inf}(S)$, $X_\varphi = \{x \in X : \mathcal{I}_x = \varphi\}$.

If $\overline{X}_\varphi = \emptyset$, φ is called an empty information, otherwise the information is non empty.

If $\overline{X}_\varphi = 1$, then φ is called selective information.

If all informations in S are selective, then S is referred to as an selective system.

If all informations in S are not empty, then S is called complete.

Definition. Let $S = \langle X, A, Q, \mathcal{I} \rangle$ be an information system. Then $\sim_a (a \in A)$ is a binary relation defined by $x \sim_a y$ iff $\mathcal{I}(x, a) = \mathcal{I}(y, a)$. And let \sim be a binary relation on X defined as:

$$x \sim y \text{ iff } \mathcal{I}_x = \mathcal{I}_y.$$

Theorem \sim_a, \sim are equivalence relations and $S = \bigcirc_{a_i \in A} \sim_{a_i}$.

where \bigcirc is a product of partitions (equivalence relations) \sim_{a_i} defined in a usual way.*

\sim_a for each $a \in A$ is an equivalence relation, because

$$1^\circ \quad x \sim_a x, \quad \text{for } \mathcal{I}_x(a) = \mathcal{I}_x(a)$$

$$2^\circ \quad x \sim_a y \Rightarrow y \sim_a x, \quad \text{because}$$

$$\mathcal{I}_x(a) = \mathcal{I}_y(a) \Rightarrow \mathcal{I}_y(a) = \mathcal{I}_x(a)$$

* Of course every subset $A' \subset A$ generates an equivalent relation

$$A' = \bigcirc_{a_i \in A'} \sim_{a_i}$$

In particular $\sim = \bigcirc_{a \in A} \sim_a$.

3^o If $x \sim_a y$ and $y \sim_a z$, then $x \sim_a z$, because if

$$\mathcal{I}_x(a) = \mathcal{I}_y(a) \quad \text{and}$$

$$\mathcal{I}_y(a) = \mathcal{I}_z(a)$$

then

$$\mathcal{I}_x(a) = \mathcal{I}_z(a)$$

Similarly, \sim is an equivalence relation;

$$1^\circ \quad x \sim x, \quad \text{because } \mathcal{I}_x = \mathcal{I}_x$$

$$2^\circ \quad \text{if } x \sim y, \text{ then } y \sim x, \quad \text{because}$$

$$\mathcal{I}_x = \mathcal{I}_y \Rightarrow \mathcal{I}_y = \mathcal{I}_x$$

$$3^\circ \quad \text{if } x \sim y \text{ and } y \sim z, \text{ then } x \sim z, \quad \text{because}$$

$$\mathcal{I}_x = \mathcal{I}_y \quad \text{and} \quad \mathcal{I}_y = \mathcal{I}_z \quad \text{implies} \quad \mathcal{I}_x = \mathcal{I}_z.$$

Finally, if

$$x \sim y, \quad \text{then } \mathcal{I}_x = \mathcal{I}_y$$

And this implies

$$\bigwedge_{a \in A} \mathcal{I}_x(a) = \mathcal{I}_y(a)$$

hence

$$\mathcal{I}_x(a_1) = \mathcal{I}_y(a_1)$$

$$\mathcal{I}_x(a_2) = \mathcal{I}_y(a_2)$$

.....

$$\mathcal{I}_x(a_n) = \mathcal{I}_y(a_n) \Rightarrow$$

$$\mathcal{I}(x \sim_{a_1} y) \wedge \mathcal{I}(x \sim_{a_2} y) \wedge \dots \wedge \mathcal{I}(x \sim_{a_n} y) \Rightarrow$$

$$\bigcirc_{a_i \in A} x \sim_{a_i} y$$

The converse implication is similar.

2. ATTRIBUTES

Definition. An information system S' is said to be included in a system S ($S' \subset S$) iff $\tilde{S}' \subset \tilde{S}$.

Definition. Two information systems S', S are said to be equivalent iff $\tilde{S}' = \tilde{S}$.

Let S', S be two information systems.

If $X = X', Q' \subset Q, A' \subset A$, and

$$S' = S /_{A \times X}$$

then we shall say that S' is weaker than S , or S is stronger than S' ($S' < S$). Of course if $S' < S$ then $S' \supset S$.

If S_1, S_2 are information systems then $S' = S_1 \cup S_2$ will be called join of S_1 and S_2 - and defined in the following manner:

$$X = X_1 \cup X_2,$$

$$A = A_1 \cup A_2,$$

$$Q = Q_1 \cup Q_2,$$

$$S = S_1 \cup S_2,$$

and

$$S /_{X_1 \times A_2} = S_1$$

$$S /_{X_2 \times A_2} = S_2$$

moreover

$$S' /_{(X_1 \cap X_2) \times (A_1 \cap A_2)} = S_2' /_{(X_1 \cap X_2) \times (A_1 \cap A_2)}.$$

Every information system $S = \langle X, A, Q, \gamma \rangle$ may be represented as join of one attribute systems $S = \bigcup_{1 \leq i \leq k} S_i$, where $Q_i =$

$$= \langle X, a_i, Q_{a_i}, S_i \rangle, \quad a_i \in A = \{a_1, \dots, a_k\}.$$

Let S_1, S_2 be information systems such that

$$X_1 = X_2, A_1 = A_2, Q_1 \neq Q_2.$$

If $a^1 \subset a^2$ for all $a \in A(A')$ we shall call S_1 finer than S_2 , or S_2 cruder than S_1 ($S_1 \sqsubset S_2$).

If $S_1 \sqsubset S_2$, then $\tilde{S}_1 \subset \tilde{S}_2$.

(\sim_1, \sim_2 are equivalence relations defined by attribute systems S_1, S_2 respectively).

Definition. Let $a, b \in A$ be two attributes in an information system S . Attribute a is said to be included in b ($a \rightarrow b$) (or b is dependent on a) iff $\tilde{a} \subset \tilde{b}$; a and b are equivalent ($a \sim b$) iff $\tilde{a} = \tilde{b}$.

If $a \rightarrow b$ then there exist a function

$$f : Q_a \rightarrow Q_b;$$

which means that the value at attribute a uniquely defines the value of attribute b .

Definition. The set A' of attributes in S is called independent in S iff for every $A' \subset A, \tilde{A'} \not\subset \tilde{A}$; if there is a subset $A' \subset A$ such that $\tilde{A'} = \tilde{A}$, then the set of attributes A' will be called dependent in S .

Definition. The smallest set $A' \subset A$ such that A' is independent in $S = \langle X, A, Q, \gamma \rangle$ will be called reduct of A , and the corresponding system, $S' = \langle X, A', Q, \gamma' \rangle$ - reduced system (γ' is the restriction of the function γ to the set A').

Theorem For every information system S, there exists an equivalent reduced information system S'.

Theorem If an information system S is complete, then S is reduced (the converse implication is not true).

Theorem If S is a reduced information system, then all its different attributes are pairwise independent (converse is not true).

Problem. Given an information system S. Find an algorithm transferring S to a reduced form.

This problem may be also stated as follows:

Let X be a set and $\Pi = \{\pi_1, \dots, \pi_n\}$ a finite family of equivalence relations on X (partitions of X) and let $\tilde{\Pi}$ be defined as

$$\tilde{\Pi} = \bigoplus_{\pi_i \in \Pi} \pi_i$$

\bigoplus - is a product of partitions.

Find the smallest subset $\Pi' \subset \Pi$ such that $\tilde{\Pi}' = \tilde{\Pi}$.

3. THE LANGUAGE OF S

With the system S we associate the language \mathcal{L}_S .

Let us first define syntax of \mathcal{L}_S .

An alphabet of \mathcal{L}_S (denoted Σ_S) is the following:

- 1° $Q \subset \Sigma_S, A \subset \Sigma_S$
- 2° $0, 1 \in \Sigma_S$
- 3° $T, F \in \Sigma_S$
- 4° $\wedge, \vee, + \in \Sigma_S$

$$5^\circ \neg, \wedge, \vee \in \sum_S$$

$$6^\circ = \in \sum_S$$

Terms \mathcal{T}_S are defined as follows

$$1^\circ 0, 1 \in \mathcal{T}_S$$

$$2^\circ A \times Q \subset \mathcal{T}_S$$

$$3^\circ \text{if } t, t' \in \mathcal{T}_S, \text{ then}$$

$$t, t \cdot t', t + t' \in \mathcal{T}_S$$

$$4^\circ \text{Nothing else is in } \mathcal{T}_S.$$

Formulae \mathcal{F}_S are the following expressions

$$1^\circ T, F \in \mathcal{F}_S$$

$$2^\circ \text{if } t, t' \in \mathcal{T}_S, \text{ then } t = t' \in \mathcal{F}_S$$

$$3^\circ \text{if } \psi, \psi' \in \mathcal{F}_S, \text{ then}$$

$$\neg \psi, \psi \vee \psi', \psi \wedge \psi' \in \mathcal{F}_S$$

$$4^\circ \text{Nothing else is in } \mathcal{F}_S.$$

Semantics of the language will consist of two parts: semantics of terms and semantics of formulae.

Semantics of terms is the function

$$\sigma_S: \mathcal{T}_S \rightarrow \mathcal{P}(X),$$

where

- 1° $\sigma_{\mathcal{F}}(0) = \emptyset, \sigma_{\mathcal{F}}(1) = X$
- 2° $\sigma_{\mathcal{F}}(a, q) = \{x \in X : x(a) = q\}$
- 3° $\sigma_{\mathcal{F}}(\sim t) = X \setminus \sigma_{\mathcal{F}}(t)$
 $\sigma_{\mathcal{F}}(t \cdot t') = \sigma_{\mathcal{F}}(t) \cap \sigma_{\mathcal{F}}(t')$

and

$$\sigma_{\mathcal{F}}(t + t') = \sigma_{\mathcal{F}}(t) \cup \sigma_{\mathcal{F}}(t').$$

Semantics of formulae is the following function

$$\sigma_{\mathcal{F}}: \mathcal{F} \rightarrow \{T, F\}.$$

here

- 1° $\sigma_{\mathcal{F}}(T) = T, \sigma_{\mathcal{F}}(F) = F$
- 2° $\sigma_{\mathcal{F}}(T=t') =$
 $\sigma_{\mathcal{F}}(t=t') = \begin{cases} T & \text{if } \sigma_{\mathcal{F}}(t) = \sigma_{\mathcal{F}}(t') \\ F & \text{if } \sigma_{\mathcal{F}}(t) \neq \sigma_{\mathcal{F}}(t') \end{cases}$
- 3° $\sigma_{\mathcal{F}}(\neg \Psi) = \begin{cases} T & \text{if } \sigma_{\mathcal{F}}(\Psi) = F \\ F & \text{if } \sigma_{\mathcal{F}}(\Psi) = T \end{cases}$
- 4° $\sigma_{\mathcal{F}}(\Psi \vee \Psi') = \sigma_{\mathcal{F}}(\Psi) \vee \sigma_{\mathcal{F}}(\Psi')$
- 5° $\sigma_{\mathcal{F}}(\Psi \wedge \Psi') = \sigma_{\mathcal{F}}(\Psi) \wedge \sigma_{\mathcal{F}}(\Psi')$

Definition. A term t is primitive if t is of the form

$$(a_1, q_1) \cdot (a_2, q_2) \cdot \dots \cdot (a_n, q_n),$$

where $A = \{a_1, \dots, a_n\}$ and $q_1, \dots, q_n \in Q$.

Definition A term t is in normal form if t is of the form

$$t_1 + t_2 + \dots + t_r,$$

with $r \geq 1$ and t_1, \dots, t_r primitive terms.

Definition Let $Y \subseteq X$. We say that Y is atomic, if there exists a primitive term t such that

$$\sigma_{\mathcal{F}}(t) = Y$$

(Note that non empty atomic sets are exactly equivalence classes of the relation \sim_S .)

Definition. Let $Y \subseteq X$. We say that Y is describable if there exists a term t such that

$$\sigma_{\mathcal{F}}(t) = Y.$$

We will use $\mathcal{T}_{\text{prim}}$ to denote the set of all primitive terms and \mathcal{T}_{nor} to denote the set of all normal terms.

Theorem. Let $t, t' \in \mathcal{T}_{\text{prim}}$ be two different primitive terms. Then

$$1^\circ \sigma_{\mathcal{F}}(t) \cap \sigma_{\mathcal{F}}(t') = \emptyset$$

$$2^\circ \bigcup_{t \in \mathcal{T}_{\text{prim}}} \sigma_{\mathcal{F}}(t) = X$$

Theorem. There exists an algorithm which for an arbitrary term t constructs a normal term t^* such that

$$\sigma_{\mathcal{Y}}(t) = \sigma_{\mathcal{Y}}(t^*)$$

Corollary. Let $Y \subseteq X$. Then Y is describable if and only if Y is a set theoretical union of atomic sets.

References

- [1] W. Marek, Z. Pawlak: Information Storage and Retrieval Systems. Mathematical Foundations. Theoretical Computer Sciences 2 (1976), pp. 331-354.
- [2] Z. Pawlak: Mathematical Foundation of Information Retrieval. CC PAS Reports, No 101 (1973).