

DAL - A LOGIC FOR DATA ANALYSIS

Luis FARINAS DEL CERRO

LSI Université Paul SABATIER
31062 TOULOUSE CEDEX

FRANCE

Ewa ORLOWSKA

Institute of Computer Science
Polish Academy of Sciences

00-901 WARSZAWA

POLAND

ABSTRACT

The aim of this paper is given some ideas how to analyze data using logic tools.

1 INTRODUCTION

The aim of this paper is to define a logic which enables us reasoning about data. Data are considered to be things or entities and properties meaningful for these entities. One of the main problems in data analysis is to induce patterns in a set of data items (Benzecri (1973)). This problem consists of two tasks :

(1) To aggregate data into sets which can be adequately characterised by means of some of the given properties.

(2) For a set of data given a priori to choose those properties from a given set of properties which are adequate for defining this set.

That is the scheme of task (1) is: from properties to sets of data items and the scheme of task (2) is : from sets of data items to properties.

The similar problems are considered in the fields of information systems (Pawlak (1981)), knowledge representation (Pawlak (1983), Konrad et al (1981), Orłowska (1983)), and pattern recognition (FU (1974)).

In our approach the formal counterparts of data are a non-empty set of objects and a family of equivalence relations on this set. Objects will be interpreted as data items and relations correspond to properties of data items. Namely, each property induces an equivalence relation such that an equivalence class of the relation consists of those objects which are the same with respect to this property.

In section 2 we give a detailed explanation of (1) and (2) on the level of semantical structure. In section 2 we introduce a syntactic structure to be used to represent the given semantical requirements. We present a language in which facts concerning tasks (1) and (2) can be formulated in a formal way. Such approach is necessary if reasoning is to be carried out by a computer. Next, we provide a deduction method for the language which enables us to prove facts concerning tasks (1) and (2).

Our approach follows the ideas and methods developed in Pawlak (1982), Orłowska and Pawlak (1981), Orłowska (1983), Harel (1978) and Mirkowska (1981).

2 DEFINABILITY OF DATA

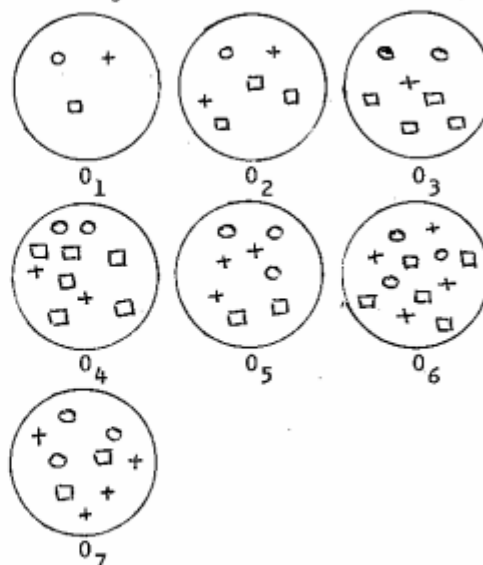
In this section we present a formal framework which enables us to express facts concerning relationships

between data items and properties of data items. We use the basic notions introduced in Konrad, Orłowska, and Pawlak (1981), Pawlak (1982), Orłowska (1983). The new idea is to consider what is called strong definability of sets of data items. It enables us to reflect an adequacy of properties for characterization of sets of data.

We consider a non-empty set OB whose elements are interpreted as data items. The elements of this set are referred to as objects. Each property meaningful for objects from set OB induces an equivalence relation in set OB . Namely, two objects are in the relation corresponding to a property if they cannot be distinguished by means of this property. Equivalence relations in set OB are referred to as indiscernibility relations. For example, if a characteristic "colour" is meaningful for elements of a set OB then we consider indiscernibility relation "to be of the same colour" which enables us to aggregate objects into classes according to their colour. These classes are equivalence classes of the indiscernibility relation in question.

Given a pair of indiscernibility relations, say R and P , we consider the intersection $R \cap P$ and the transitive closure $R \circ P$ of these relations. Clearly $R \cap P$ and $R \circ P$ are equivalence relations. Equivalence classes of $R \cap P$ and $R \circ P$ are obtained from equivalence classes of R and P by making their intersections and unions, respectively. Relations $R \cap P$ and $R \circ P$ can be considered to be indiscernibility relations corresponding to properties which are composed in some way from properties related to R and P . Let us consider a simple example.

Example 2.1. Assume that we are given seven objects :



The natural properties meaningful for these objects are : number of circles, number of crosses, number of squares. Indiscernibility relations R_o , R_+ and R_{\square} corresponding to these properties provide the following classes :

$$R_o \{0_1, 0_2\} \{0_3, 0_4\} \{0_5, 0_6, 0_7\}$$

$$R_+ \{0_1, 0_3\} \{0_2, 0_4\} \{0_6, 0_7\} \{0_5\}$$

$$R_{\square} \{0_1\} \{0_2\} \{0_3\} \{0_4\} \{0_5, 0_7\} \{0_6\}$$

The classes determined by relation $R_o \cap R_+$ are as follows :

$$R_o \cap R_+ \{0_1\} \{0_2\} \{0_3\} \{0_4\} \{0_5\} \{0_6, 0_7\}$$

The classes reflect the pattern corresponding to property "to have the same number of circles and the same number of crosses".

The classes corresponding to relation $R_o \circ R_+$ are as follows :

$$R_o \circ R_+ \{0_1, 0_2, 0_3, 0_4\} \{0_5, 0_6, 0_7\}$$

The property corresponding to this pattern can be expressed as "number of crosses and number of circles less than 3".

In general, a property corresponding to the intersection of indiscerni-

bility relations R and P can be defined as a conjunction of properties corresponding to R and P . In the case of $R \oplus P$ there is no such straightforward relationship between the underlying properties. However the property corresponding to $R \oplus P$ can be defined as a certain relation between properties corresponding to R and P .

Given an indiscernibility relation in a set OB of data items, we consider the problem of definability of subsets of set OB in terms of the property corresponding to this relation. Let X and R be a subset of set OB and an indiscernibility relation in set OB , respectively.

A lower approximation $\underline{R}X$ of set X with respect to relation R is a union of those equivalence classes of R which are included in X .

An upper approximation $\overline{R}X$ of set X with respect to relation R is a union of those equivalence classes of R which have an element in common with X .

In terms of approximations we can define positive, negative, and borderline instances of a set, namely :

- $\underline{R}X$ set of positive instances of X with respect to R
- $OB - \overline{R}X$ set of negative instances of X with respect to R
- $\overline{R}X - \underline{R}X$ set of borderline instances of X with respect to R .

By using the notions given above we can express facts concerning the kind of data analysis in which we are interested in the structure within a given a priori set of data. In this kind of task we assume that we are given a subset X of data items determined by means of a certain external condition and we are interested in establishing a relationship between X and the structuring of data provided by indis-

cernibility relations.

We say that set X is definable with respect to indiscernibility R iff $\underline{R}X = X = \overline{R}X$.

In other words a set is definable with respect to indiscernibility relation iff it can be covered by equivalence classes of this relation. This means that a pattern provided by X can be expressed by means of a property corresponding to this indiscernibility. If a set X is not definable with respect to an indiscernibility R then we can express the pattern given by X with a certain degree of inexactness. Set of positive instances of X with respect to R is the greatest definable set included in R which represents the pattern corresponding to X , set of negative instances of X is the greatest definable set whose elements do not obey this pattern, set of borderline instances consists of elements for which we cannot decide whether they obey the pattern or not.

We say that set X is strongly definable with respect to indiscernibility R iff $\underline{R}X = X = \overline{R}X$ and X is an equivalence class of R .

We make distinction between a possibility of covering a set by one equivalence class and by more than one class to reflect an adequacy or properties for characterization of sets of data. If a set is strongly definable with respect to a relation then we will consider a property corresponding to this relation to be adequate for describing the set.

Example 2.2. Let us consider the set $OB = \{0_1, 0_2, \dots, 0_7\}$ given in example 2.1., and indiscernibility relations R_0, R_+ , and R_{\square} . Set $X = \{0_1, 0_2, 0_3, 0_4\}$ is definable with respect to R_0, R_+ , and

R_{\square} , but it is not strongly definable with respect to these relations. However X is strongly definable with respect to relation $R_0 \cup R_+$. Set $Y = \{0_1, \dots, 0_2, 0_7\}$ is not definable neither by R_0, R_+ , and R_{\square} , nor by $R_0 \cap R_+$, and $R_0 \cap R_{\square}$. Set Y is definable with respect to $R_+ \cap R_{\square}$, but none of the relations obtained from R_0, R_+ , and R_{\square} by using operations \cup and \cap is sufficient to provide the strong definability of Y . The approximations of set Y with respect to some of the indiscernibility relations in question are given below.

- $\underline{R_0} Y = \{0_1, 0_2\}$
- $\overline{R_0} Y = \{0_1, 0_2, 0_5, 0_6, 0_7\}$ positive instances of Y with respect to $R_0 : 0_1, 0_2$. Negative instances of Y with respect to $R_0 : 0_3, 0_4$. Borderline instances of Y with respect to $R_0 : 0_5, 0_6, 0_7$
- $\underline{R_0 \cap R_+} Y = \{0_1, 0_2\}$
- $\overline{R_0 \cap R_+} Y = \{0_1, 0_2, 0_6, 0_7\}$

In terms of the notions defined above we can discuss the ability or properties to define sets of data. We can consider a certain property P to be better than a property P' for characterization of a set X of data iff the approximations of X with respect to the indiscernibility determined by P are closer to X (with respect to inclusion) than the approximations with respect to the indiscernibility corresponding to P' .

In the following we list some properties of approximations.

$$\begin{aligned} \underline{RX} &\subseteq X \\ \underline{RRX} &= \underline{RX} \\ \underline{R(X \cap Y)} &= \underline{RX} \cap \underline{RY} \\ \underline{RX \cup RY} &\subseteq \underline{R(X \cup Y)} \\ \text{If } X &\subseteq Y \text{ then } \underline{RX} \subseteq \underline{RY} \\ \underline{RX \cup SX} &\subseteq \underline{R(X \cup SX)} \\ \underline{X} &\subseteq \underline{RX} \\ \underline{RRX} &= \underline{RX} \end{aligned}$$

$$\begin{aligned} \overline{R(X \cup Y)} &= \overline{RX} \cup \overline{RY} \\ \overline{R(X \cap Y)} &= \overline{RX} \cap \overline{RY} \\ \text{If } X &\subseteq Y \text{ then } \overline{RX} \subseteq \overline{RY} \\ \underline{R} \cup \underline{S} X &\subseteq \underline{RX} \cap \underline{SX} \end{aligned}$$

In the next section we present a formal language in which we can express facts concerning definability of data discussed in the present section. The set-theoretical notions introduced here provide a basis for semantics of the language.

3 THE LANGUAGE OF LOGIC DAL

Expressions of the language of logic DAL are built from the symbols of the following pairwise disjoint sets :

- VAR propositional variables
- VARREL relational variables
- $\neg, \wedge, \vee, \leftrightarrow, \rightarrow$ classical propositional operations of negation, conjunction, implication and equivalence, respectively.
- \cup, \cap binary operations on relations
- [], < > unary modal propositional operations
- () brackets

We assume that sets VAR and VARREL are non-empty, at most denumerable sets.

Set EREL of relational expressions is the least set satisfying the following conditions :

$$\begin{aligned} \text{VARREL} &\subseteq \text{EREL} \\ R, S \in \text{EREL} &\text{ implies } R \cup S, R \cap S \in \text{EREL} \end{aligned}$$

Relational variables are intended to represent indiscernibility relations, and operations \cap and \cup will be interpreted as the intersection and the transitive closure of the union of relations.

Set FOR of all formulas of the language is the least set satisfying the following conditions :

$$\begin{aligned} \text{VAR} &\subseteq \text{FOR} \\ A, B \in \text{FOR} &\text{ implies } \neg A, A \vee B, A \wedge B, \end{aligned}$$

$A \rightarrow B, A \leftrightarrow B \in \text{FOR}$
 $A \in \text{FOR}$ and $R \in \text{EREL}$ imply $[R]A,$
 $\langle R \rangle A \in \text{FOR}.$

Formulas are intended to represent sets of data items. In particular formulas built by using modal operations correspond to approximations of sets. Since in the language of DAL we allow compound relational expressions, we can express relationships between approximations with respect to various properties and we can explicitly describe indiscernibility relations corresponding to these properties.

4 SEMANTICS OF THE LANGUAGE OF LOGIC DAL

To define meaning of formulas of logic DAL we should fix a set OB of data items and a family of equivalence relations in set OB corresponding to properties of these data. To be more formal, we define the notions of model and satisfiability of the formulas in a model. By a model we mean a triple : $M = (OB, \{\rho_R\}_{R \in \text{EREL}}, m)$, where OB is a non-empty set; for any $R \in \text{EREL}$ ρ_R is an equivalence relation in set OB such that $\rho_R \cap \rho_S$ is the greatest equivalence relation in set OB included both in ρ_R and ρ_S ; $\rho_R \oplus \rho_S$ is the least equivalence relation including both ρ_R and ρ_S ; $m : \text{VAR} \cup \text{VARREL} \rightarrow \mathcal{P}(OB) \cup \{\rho_R\}_{R \in \text{EREL}}$ is a meaning function such that :

- . $m(p) \subseteq OB$ for $p \in \text{VAR}$
- .. $m(R) = \rho_R$ for $R \in \text{EREL}$

Given a model M we say that a formula A is satisfied by an object o in model M ($M, o \text{ sat } A$) iff the following conditions are satisfied :

$M, o \text{ sat } P$ iff $o \in m(p)$ for $p \in \text{VAR}$
 $M, o \text{ sat } \neg A$ iff not $M, o \text{ sat } A$
 $M, o \text{ sat } A \vee B$ iff $M, o \text{ sat } A$ or
 $M, o \text{ sat } B$

$M, o \text{ sat } A \wedge B$ iff $M, o \text{ sat } A$ and
 $M, o \text{ sat } B$

$M, o \text{ sat } A \rightarrow B$ iff $M, o \text{ sat } (\neg A \vee B)$
 $M, o \text{ sat } A \leftrightarrow B$ iff $M, o \text{ sat}$
 $(A \rightarrow B) \wedge (B \rightarrow A)$

$M, o \text{ sat } [R] A$ iff for all $o' \in OB$
if $(o, o') \in \rho_R$ then $M, o' \text{ sat } A$
 $M, o \text{ sat } \langle R \rangle A$ iff there is an
 $o' \in OB$ and that $(o, o') \in \rho_R$ and
 $M, o' \text{ sat } A.$

Given a model M , to each formula A of the language we assign a set of objects called an extension of A in model M ($\text{ext}_M A$) :

$\text{ext}_M A = \{o \in OB : M, o \text{ sat } A\}$

The immediate consequences of this definition are the following :

Fact 4.1.

- (a) $\text{ext}_M p = m(p)$ for $p \in \text{VAR}$
- (b) $\text{ext}_M \neg A = -\text{ext}_M A$
- (c) $\text{ext}_M (A \vee B) = \text{ext}_M A \cup \text{ext}_M B$
- (d) $\text{ext}_M (A \wedge B) = \text{ext}_M A \cap \text{ext}_M B$
- (e) $\text{ext}_M (A \rightarrow B) = -\text{ext}_M A \cup \text{ext}_M B$
- (f) $\text{ext}_M (A \leftrightarrow B) = (\text{ext}_M A \cap \text{ext}_M B) \cup (-\text{ext}_M A \cap -\text{ext}_M B)$
- (g) $\text{ext}_M [R] A = \rho_R \text{ ext}_M A$
- (h) $\text{ext}_M \langle R \rangle A = \overline{\rho_R} \text{ ext}_M A$

Hence the classical propositional operations are interpreted as set-theoretical operations and the modal operations correspond to the operations of lower and upper approximation.

We admit the usual notions of truth and validity of formulas. A formula A is true in a model M ($\vDash_M A$) iff $\text{ext}_M A = OB$. A formula A is valid in logic DAL ($\vDash A$) iff A is true in every model for DAL. A formula A is a semantical consequence of a set Γ of formulas ($\Gamma \vDash A$) iff for any model M formula A is true in M whenever all formulas from Γ are true in M . A formula A is satisfiable iff $M, o \text{ sat } A$ for some model M and object o . A set Γ of formulas is satisfied in a model M by an object

$o (M, o \text{ sat } \Gamma)$ iff $M, o \text{ sat } A$ for all $A \in \Gamma$. A set Γ is satisfiable iff $M, o \text{ sat } \Gamma$ for some model M and object o .

Given a model $M = (OB, \{\rho_R\}_{R \in \text{EREL}}, m)$, meaning function m provides a family of sets of data items which we are interested in. Next, we consider compound sets, which are expressed by means of formulas obtained from propositional variables by performing classical propositional operations. We can also express approximations of these sets with respect to relations admitted in the model. As a consequence we can discuss on a formal level definability of data in terms of properties related to these relations.

In the following we show how we can express facts concerning sets of data by means of formulas of DAL. As usually, we can express inclusion and equality of sets.

Fact 4.2.

(a) $\overline{M} A \rightarrow B$ iff $\text{ext}_M A \subseteq \text{ext}_M B$

(b) $\overline{M} A \leftrightarrow B$ iff $\text{ext}_M A = \text{ext}_M B$

definability and strong definability of sets of data can be expressed as follows.

Fact 4.3.

(a) $\overline{M} \langle R \rangle A \rightarrow [R]A$ iff $\text{ext}_M A$ is definable with respect to relation ρ_R

(b) $\overline{M} (B \rightarrow [R]A) \rightarrow (([R]B \rightarrow A \wedge \neg A) \vee ([R]B \rightarrow [R]A))$ for every

formula B iff $\text{ext}_M [R]A$ is an equivalence class of relation ρ_R or the empty set.

The formula in condition (a) assures that the upper approximation of set $\text{ext}_M A$ is included in the lower approximation and hence by the definition of definability set $\text{ext}_M A$ is definable with respect to ρ_R . The family of formulas in condition (b) assures that the

lower approximation of any set of data properly included in the lower approximation of set $\text{ext}_M A$ equals the empty set. This means that the lower approximation of $\text{ext}_M A$ consists of exactly one equivalence class of relation ρ_R . Hence the formulas from condition (a) and the family of formulas from condition (b) assure the strong definability of set $\text{ext}_M A$.

Fact 4.4.

$\overline{M} [R]A \rightarrow [P]A$ and $\overline{M} \langle P \rangle A \rightarrow \langle R \rangle A$ iff $\text{ext}_M A$ is characterised better by ρ_P than by ρ_R .

The above formulas express inclusions of approximations of set $\text{ext}_M A$ with respect to relations ρ_P and ρ_R . The lower approximation with respect to ρ_P is greater than the lower approximation with respect to ρ_R , and the upper approximation with respect to ρ_P is smaller than the upper approximation with respect to ρ_R . This means that the approximations with respect to ρ_P are closer to set $\text{ext}_M A$ than the approximations with respect to ρ_R .

Fact 4.5.

$\overline{M} [R]A \rightarrow [P]A$ and $\overline{M} \langle P \rangle A \rightarrow \langle R \rangle A$ for every formula A iff $\rho_P \leq \rho_R$

The condition says that inclusion of indiscernibility relations can be expressed by a family of formulas which assure that for any set of data lower approximations with respect to ρ_P are greater than the lower approximations with respect to ρ_R , and the upper approximations with respect to ρ_P are smaller than the upper approximations with respect to ρ_R .

We conclude, that in the language of logic DAL we can express various kinds of information: facts concerning sets of data items, facts concerning indiscernibility relations corresponding

to properties of data items, and relationships between data items and their properties, especially those concerning definability.

5 AXIOMATIZATION

In this section we present a deductive system for the language of DAL. We admit the following schemes of axioms and inference rules. Let R , S , and P denote arbitrary relational expressions and let A , B denote formulas.

Axioms of DAL

- A1. All formulas having the form of a tautology of the classical propositional logic
 A2. $[R] (A \rightarrow B) \rightarrow ([R]A \rightarrow [R]B)$
 A3. $[R]A \rightarrow A$
 A4. $\langle R \rangle A \rightarrow [R] \langle R \rangle A$
 A5. $[R \cup S]A \rightarrow [R]A \wedge [S]A$
 A6. $(([P]A \rightarrow [R]A) \wedge ([P]A \rightarrow [S]A)) \rightarrow ([P]A \rightarrow [R \cup S]A)$
 A7. $[R]A \vee [S]A \rightarrow [R \cap S]A$
 A8. $(([R]A \rightarrow [P]A) \wedge ([S]A \rightarrow [P]A)) \rightarrow ([R \cap S]A \rightarrow [P]A)$

Rules of inference

- | | |
|----------------------|---------------|
| $A, A \rightarrow B$ | modus ponens |
| B | |
| A | necessitation |
| $[R]A$ | |

For a fixed relation R axioms A1, ..., A4 and the rules of inference correspond to the axiomatization of modal logic S5. Axioms A5, A6 provide the definition of operation \cup and axioms A7, A8 give the definition of operation \cap .

In the usual way we define the notions of proof and theorem. A proof of a formula A from a set Γ of formulas is a finite sequence of formulas each of which is either an axiom or an element of set Γ or else is obtainable from earlier formulas by a rule of inference. A formula A is derivable from

a set Γ ($\Gamma \vdash A$) iff it has a proof from set Γ . A formula A is a theorem of DAL ($\vdash A$) iff it is derivable merely from axioms. A set Γ of formulas is consistent if the formula of the form $A \wedge \neg A$ is not derivable from Γ .

It is easy to see that the axioms are valid and the rules preserve validity. Hence the following theorem holds.

Fact 5.1. (Soundness theorem)

- (a) $\vdash A$ implies $\vDash A$
 (b) $\Gamma \vdash A$ implies $\Gamma \vDash A$
 (c) Γ satisfiable implies Γ consistent.

In the following we list some theorems of logic DAL.

Fact 5.2.

- (a) $\vdash [R \cup R]A \leftrightarrow [R]A$
 (b) $\vdash [R \cap R]A \leftrightarrow [R]A$
 (c) $\vdash [(R \cup S) \cup P]A \leftrightarrow [R \cup (S \cup P)]A$
 (d) $\vdash [(R \cap S) \cap P]A \leftrightarrow [R \cap (S \cap P)]A$
 (e) $\vdash [(R \cup S) \cap P]A \rightarrow [(R \cap P) \cup (S \cap P)]A$
 (f) $\vdash (R \cup P) \cap (S \cup P)A \rightarrow [(R \cap S) \cup P]A$
 (g) $\vdash [R \cup S]A \rightarrow [R][S]A$
 (h) $\vdash [R][S]A \rightarrow [R \cap S]A$

Fact 5.3. (Completeness theorem)

- (a) $\vDash A$ implies $\vdash A$
 (b) $\Gamma \vDash A$ implies $\Gamma \vdash A$
 (c) Γ consistent implies Γ satisfiable.

6 CONCLUSIONS

In this paper we have given some ideas how to analyse data using logic tools. We followed the classical framework of logic programming, namely we defined a logic language for expressing data analysis problems and we developed a deductive system for the language to use proof procedures to obtain solutions to these problems.

Data analysis has been understood

as a process of obtaining patterns in a set of data items. We considered two main tasks involved in data analysis :

- . to aggregate data into sets according to their properties
- . to define properties adequate for characterisation of sets of data.

We defined formal counterparts of sets of data and properties. Namely, sets of data are defined by means of formulas of the language of logic DAL and properties are defined by means of relational expressions of the language. We presented the deductive system for logic DAL. We have given a particular interest to the notion of strong definability of sets of data which enables us to establish properties which adequately characterize these data.

The results of the present paper can be extended to languages with the other operations on relations e.g. with the composition of relations, or to relations which are not necessarily equivalence relations e.g. tolerance relations (reflexive and symmetric). We expect that a mechanical proof procedure can be defined for logic DAL by u-

Konrad, E., Orłowska, E., Pawlak, Z. Knowledge representation systems. ICSPAS Report 433, 1981.

Mirkowska, G. PAL - Propositional algorithmic logic, logic of programs. Lecture Notes in Computer Science, 31-101, Springer-Verlag, 1981.

Orłowska, E. Representation of vague information. ICSPAS Reports 503, 1983.

Orłowska, E., Pawlak, Z. Expressive power of knowledge representation systems. ICSPAS Reports 432, accepted for publication in the Journal of Man Machine Studies, 1981.

Pawlak, Z. Information systems-theoretical foundations. Information Systems 6, 205-218, 1981.

Pawlak, Z. Rough sets. International Journal of Computer and Information Sciences, 11, 341-356, 1982.

Pawlak, Z. Rough classification. ICSPAS Reports 506, accepted for publication in the Journal of Man-Machine Studies, 1983.