# A Method for Investigating Links between Discrete Data Features in Knowledge Bases in the Form of Predicate Equations

Oleksandr Karataiev[1], Dmytro Sitnikov[1] and Nataliia Sharonova[2]

[1] *Kharkiv National University of Radio Electronics, Nauky ave, 14, Kharkiv, 61166, Ukraine*
[2] *National Technical University "Kharkiv Polytechnic Institute", Kyrpychova str., 2, Kharkiv, 61000, Ukraine*

### Abstract

Logical classification methods normally involve the compilation and solution of logical equations with variables that take values of 1 and 0, depending on whether the given object has a certain property or not. The solution of such equations makes it possible either to identify an object by the available sets of values of attribute variables, or to establish unknown properties of a given object. A natural generalization of the Boolean algebra equations is equations of the finite predicate algebra, which makes it possible to operate with arbitrary attribute variables defined on different finite sets. The use of such equations for constructing logical conclusions in knowledge bases allows expanding the capabilities of Boolean logical methods for object recognition and classification. When classifying objects, one deals with sets of features, selecting some values of which it is possible to identify whether the object under consideration belongs to a certain class. In this paper a method for investigating links between discrete object features is presented. Also, different types of predicate equations are considered. When analyzing links between salient data features, we often encounter quite complicated systems of logic equations that, nevertheless, can be simplified owing to their specific properties. A real-world medical example has been considered to demonstrate the procedure of eliminating non-salient features.

### Keywords

Knowledge representation, classification of discrete objects, feature selection, systems of predicate equations, variables exclusion

## 1. Introduction

Knowledge representation and interpretation plays an important part in various fields of computer science. To formalize information about objects and processes in knowledge bases, various methods of discrete mathematics are used. In cases where information about objects and processes, represented by discrete information features, has a rather complex logical structure, various methods and models of discrete mathematics, including logical equations with Boolean variables, are used for its formal presentation.

Logical classification methods normally involve the compilation and solution of logical equations with variables that take values of 1 and 0, depending on whether the given object has a certain property or not. The solution of such equations makes it possible either to identify an object by the available sets of values of attribute variables, or to establish unknown properties of a given object. A natural generalization of the Boolean algebra equations is equations of the finite predicate algebra, which makes it possible to operate with arbitrary attribute variables defined on different finite sets.

The use of such equations for constructing logical conclusions in knowledge bases allows expanding the capabilities of Boolean logical methods for object recognition and classification. When classifying objects, one deals with sets of features, selecting some values of which it is possible to identify whether or not the object under consideration belongs to a certain class. In this paper a method for investigating links between discrete object features is presented. Also, different types of predicate equations are considered. When analyzing links between salient data features, we often encounter quite complicated systems of logic equations that, nevertheless, can be simplified owing to their specific properties. Excluding extra variables with the help of the quantifiers leads to the simplification of the original system of predicate equations. A real-world medical example has been considered to demonstrate the procedure of eliminating non-salient features.

## 2. Related works

Many practical problems lead to the necessity of using logic classification methods. For example, in [1] binary feature vectors are classified. The proposed method can be used in a variety of classification problems in different industries. The process of classification is reduced to investigating logical-dynamic systems depending on some initial states. In [2] corrective functions for object recognition logical methods are constructed. An interesting algorithm for logic classification with the help of correcting functions has been proposed. In [3] logical data classification has been used for the analysis of hyperspectral data. Some combinatorial issues of logic classification have been highlighted in [4]. The authors have applied their research to Alzheimer's Disease Proteomics Expression classification. In [5] logical algorithmic methods for building decision trees have been considered. A logic-based classification method for text recognition has been proposed in [6]. In [7] logic machine learning approach has been suggested. The authors have developed a special logic classifier. Logic-based design of a strong classifier with the help of weak classifiers has been considered in [8]. A combination of various methods including logic classification has been tested in [9]. A method for the classification of text messages based on logical extraction has been considered in [10]. In [11] finite predicate networks have been investigated for radar detection. A new approach to logic classification and recognition has been suggested in [12]. Subsystems of Boolean equation systems have been studied in [13]. A method for distributed solving logic equations has been presented in [14]. In [15] some specific types of logic functions are considered. The focus if on symmetric functions that are widely used for classification purposes. Entropy issues in Boolean networks have been highlighted in [16]. Logic-predicate networks have been investigated in [17]. Using predicates for classifying access issues on the Internet of things has been considered in [18]. Many works devoted to classification often use Boolean algebra, fuzzy logic or neural networks as main mathematical tools [19]. Logic approaches to fact-based analysis have been highlighted in [20].

## 3. Methodology for analyzing logic links between salient data features in systems of predicate equations

A universal way to solve systems of equations of the algebra of finite predicates is to reduce the predicate given by the system of equations and initial conditions to a perfect disjunctive normal form. However, such a procedure involves enumeration of many intermediate solutions, and its practical implementation requires a significant amount of computer time. For some types of predicate equations, taking into account the peculiarities of their structure, it is possible to develop simpler algorithms for solving them.

In many practical tasks related to the semantic processing of medical data, natural language information, customer data, there is no need to obtain all sets of values of semantic features, but it is required to obtain one or more sets of values of features (target variables) that are of interest to the user. It is often necessary to find the values of target variables under given initial conditions, which are a fixed set of values of other features. When solving such problems, other variables that are not

included in the initial conditions and are not target variables are excluded from the equation by linking them with existential quantifiers.

Unlike Boolean variables, predicate variables provide more flexibility in discovering necessary features. For example, let us consider the following dependencies:

$$y^{a_1} \rightarrow x_1^{b_1} \lor x_1^{b_2},$$
$$y^{a_2} \rightarrow x_1^{b_3} \lor x_1^{b_4},$$
$$y^{a_3} \rightarrow x_1^{b_5} \lor x_1^{b_6},$$

where domains for $y$ and $x_1$ are $\{a_1, a_2, a_3\}$ and $\{b_1, b_2, b_3, b_4, b_5, b_6\}$ correspondingly. If $y = a_1$, then $x_1 = b_1$ or $x_1 = b_2$. On the other hand, it follows from the first expression, that

$$\neg\left(x_1^{b_1} \lor x_1^{b_2}\right) \rightarrow \neg y^{a_1},$$

which means

$$x_1^{b_3} \lor x_1^{b_4} \lor x_1^{b_5} \lor x_1^{b_6} \rightarrow y^{a_2} \lor y^{a_3}.$$

Thus, if $x_1$ takes on a value from the set $\{b_3, b_4, b_5, b_6\}$, the object property $y$ takes on values either $a_2$ or $a_3$.

Finite predicates algebra gives us an opportunity to interpret knowledge in a strict mathematical form, where different features and their values are connected with the help of Boolean and predicate operations. The classical form of a logic equation with finite predicates is as follows:

$$f(x_1, x_2, \dots, x_n) = 1,$$

where each variable takes on values from a finite set of elements, and in the general case these domains can be different. In practice, we can encounter problems with many equations. In this case a problem can be resolved by solving a system of equations:

$$f_1(x_1, x_2, \dots, x_n) = 1,$$
$$f_2(x_1, x_2, \dots, x_n) = 1,$$
$$\dots$$
$$f_m(x_1, x_2, \dots, x_n) = 1.$$

If needed, this system can be rewritten as a conjunction of the above equations and represented in the form of a single equation:

$$f_1(x_1, x_2, \dots, x_n) \land f_2(x_1, x_2, \dots, x_n) \land \dots \land f_m(x_1, x_2, \dots, x_n) = 1.$$

For such equations we can define some problems that can be resolved:
1. Find all possible sets of variable values that satisfy these equations. This problem is obviously difficult as there is exponential growth in calculations.
2. Determine whether the system has a solution
3. Determine whether it has a single solution
4. Find some important combinations of variable values that satisfy the system
5. Solve the system under some initial conditions.

Consider the following system of predicate equations:

$$y^{a_1} \rightarrow g_1(x_1, x_2, \dots, x_n),$$
$$y^{a_2} \rightarrow g_2(x_1, x_2, \dots, x_n),$$
$$\dots$$
$$y^{a_m} \rightarrow g_m(x_1, x_2, \dots, x_n).$$

It means that when the feature $y$ takes on a value $a_i$, a set of values for the variables $x_1, x_2, \ldots, x_n$ should satisfy the equation

$$g_i(x_1, x_2, \ldots, x_n) = 1, \tag{1}$$

which means that if an object possesses the property $a_i$ then the features $x_1, x_2, \ldots, x_n$ should satisfy the above equation.

Generally speaking, the converse is not necessarily true. If $g_i(x_1, x_2, \ldots, x_n) = 1$, $y$ does not necessarily take on the value $a_i$. Consider a stronger dependence:

$$y^{a_1} = g_1(x_1, x_2, \ldots, x_n),$$
$$y^{a_2} = g_2(x_1, x_2, \ldots, x_n),$$
$$\ldots$$
$$y^{a_m} = g_m(x_1, x_2, \ldots, x_n).$$

In this case any set of values of the features $x_1, x_2, \ldots, x_n$ either confirms the fact that $y$ equals $a_i$ or not (belongs to the corresponding class or not). If the object has the property $a_i$, its features should satisfy (1). Thus, we can classify the feature $y$ by values of the features $x_1, x_2, \ldots, x_n$. Also, it can be easily shown that the conjunction of any two different functions $g_i$ and $g_j$ is equal to zero. It follows from the basic properties of recognition predicates. It can be shown that the above system is equivalent to the following equation:

$$y^{a_1} g_1(x_1, x_2, \ldots, x_n) \vee y^{a_2} g_2(x_1, x_2, \ldots, x_n) \vee \ldots \vee y^{a_m} g_m(x_1, x_2, \ldots, x_n) = 1. \tag{2}$$

Logic methods for object classification are applied for solving practical problems from a variety of fields: biology, physics, meteorology etc. Their specifics can be discovered at the stage of building a mathematical model including data features. Normally propositional logic is used for this purpose. We suggest an approach based on finite predicate algebra. Let us build a general form of such tasks based on predicate equations.

Let feature variables $y_1, y_2, \ldots, y_l$ denote some properties of objects, for example, a disease like flue. Each variable takes on its values from its domain. Unlike Boolean variables, predicate variables can take on values from different domains.

Let discrete variables $x_1, x_2, \ldots, x_n$ be features by sets of which we can determine which values the property variables can take on. Properties and features can be connected in the form of some complicated logic dependencies that can be represented as a predicate equation:

$$P(y_1, y_2, \ldots, y_l; x_1, x_2, \ldots, x_n) = 1, \tag{3}$$

where P is a finite predicate.

To classify an object under consideration means to determine based on this predicate equation and experimental data on the features $x_1, x_2, \ldots, x_n$, which properties (values of the features $y_1, y_2, \ldots, y_l$) this object possesses, and which properties are not satisfied. Each elementary conjunction, for example,

$$x_1^{a_{11}} x_2^{a_{12}} \ldots x_n^{a_{1n}},$$
$$x_1^{a_{21}} x_2^{a_{22}} \ldots x_n^{a_{2n}},$$
$$\ldots$$
$$x_1^{a_{m1}} x_2^{a_{m2}} \ldots x_n^{mn}$$

characterizes an object class. Then, based on the a priori dependence (3) and experimental data on the features $x_1, x_2, \ldots, x_n$, it is possible to determine to which class the given object belongs. As can be seen from the above considerations, the values of features are grouped into a matrix.

Suppose that as an experiment outcome we have obtained some data related to values of the features $x_1, x_2, \ldots, x_n$ that describe the object being classified and composed the following predicate equation describing links between them:

$$g(\,x_1, x_2, \dots, x_n) = 1.$$

The problem of object classification can be formalized as solving the following predicate equation by finding an unknown predicate $f$:

$$g(\,x_1, x_2, \dots, x_n) \rightarrow f(y_1, y_2, \dots, y_l).$$

By solving this functional equation, it is possible to determine feature values $x_1, x_2, \dots, x_n$ that characterize the objects $y_1, y_2, \dots, y_l$.

In studies related to logical inferences in knowledge bases, questions arise in determining the tightness of the links between the features of these objects, as well as questions of their materiality and insignificance. Apparently, we can consider the formal relationship between features to be stronger, the fewer sets of values of these variables satisfy the equation. In this case, if any sets of values of these variables satisfy the original equation, we can assume that there is no connection between these variables.

In addition, when solving practical problems, the following questions arise:
1. How will the specific values of this feature, substituted into the logical equation, affect the links between the other features?
2. How strong is the logical relationship between two (or more) given features?

To answer the first question, it seems natural to single out those predicates (and, accordingly, equations) that, when a certain attribute value is substituted, are transformed into predicates that give a stronger connection between variables, as well as such predicates, substitution into which this value leads to a weakening of the logical connection between signs.

To get an answer to the second question, it is necessary to exclude from the original equation with the help of the existence or universal quantifier all variables except those under consideration and study the resulting equation with a smaller number of variables, which describes all admissible sets of values of the features under study.

Let us consider the procedure of feature selection, where the number of features can be reduced. Here we can encounter the following problems:

We may need to find some sets of feature values that interest us where there is at least one value of non-salient features such that there exists at least one set of values of salient features. In this case we apply an existence quantifier to the set of non-salient values.

We may need to find some sets of feature values, where for any set of non-salient features there exists at least one solution of the equation. In this case we apply a universal quantifier to non-salient variables.

We may need to find some sets of feature values that satisfy the equation under the condition that non-salient features take on some specific values.

Let predicate P depend on the variables $x, y, \dots, z$. Define the substitution operator a(P) (a belongs to the domain of the definition of the variable x ) acting on the predicate P as follows:
$$a(P(x, y, \dots, z) = P(a, y, \dots, z).$$

Let's call the substitution operator restrictive if the following condition is met
$$P(a, y, \dots, z) \rightarrow P(x, y, \dots, z)$$
for all $x, y, \dots, z$.

Call the subsitution operator distributive if the condition is met

$$P(a, y, \dots, z) \leftarrow P(x, y, \dots, z)$$
for all $x, y, \dots, z$.

When interpreting knowledge represented by this implication, we can say that estrictive operators reinforce the logical relationship between discrete features, distributing substitution operators weaken this relationship, shifting the relationship between features in an arbitrary way.

Consider the predicate P as follows:

$$P(x, y, \ldots, z) = x^{a_1} P_1(y, \ldots, z) \vee x^{a_2} P_2(y, \ldots, z) \vee \ldots$$
$$\vee\, x^{a_n} P_n(y, \ldots, z).$$

Then

$$a_1(P) = P_1(y, \ldots, z) = x^{a_1} P_1(y, \ldots, z) \vee x^{a_2} P_1(y, \ldots, z) \vee \ldots$$
$$\vee\, x^{a_n} P_1(y, \ldots, z).$$

It is obvious that the predicate $a_1(P)$ will be contracting, if $P_1 \to P_i \forall i = 1,2,\ldots,n$.
The operator $a_1(P)$ will be distributing, if $P_1 \leftarrow P_i \forall i = 1,2,\ldots,n$.
Let us consider examples of the application of the operator $a_1$ to the predicate $P(x,y)$, where the variables $x, y$ and $z$ have the domains $\{a_1, a_2\}$, $\{b_1, b_2\}$ и $\{c_1, c_2\}$ correspondingly.
Let

$$P = x^{a_1} y^{b_1} z^{c_1} \vee x^{a_2} y^{b_1} z^{c_2} \vee x^{a_2} y^{b_1} z^{c_1}.$$

Then

$$a_1(P) = y^{b_1} z^{c_1} = (x^{a_1} \vee x^{a_2}) \& y^{b_1} z^{c_1} =$$
$$= x^{a_1} y^{b_1} z^{c_1} \vee x^{a_2} y^{b_1} z^{c_1}.$$

Except for the disjuncts that the predicate $a_1(P)$ contains $P$ includes one more disjunct $x^{a_2} y^{b_1} c^{c_1}$, i.e. the operator $a_1$ is a restricting one for the predicate $P$. According to the introduced definitions, in the given example $P_1 = y^{b_1} z^{c_1}$, $P_2 = y^{b_1} z^{c_2} \vee y^{b_1} z_1^{c_1}$ . It is obvious here that $P_1 \to P_2$. Consider now the predicate

$$P = x^{a_1} y^{b_1} z^{c_1} \vee x^{a_1} y^{b_1} z^{c_2} \vee x^{a_2} y^{b_1} z^{c_1},$$
$$a_1(P) = y^{b_1} z^{c_1} \vee y^{b_1} z^{c_2} = (x^{a_1} \vee x^{a_2}) \&$$
$$\& \left( y^{b_1} z^{c_1} \vee y^{b_1} z^{c_2} \right) = x^{a_1} y_1 b_1 c_1 \vee x^{a_2} y^{b_1} z^{c_2} \vee$$
$$\vee\, x^{a_2} y^{b_1} z^{c_1} \vee x x^{a_2} b_1 b^{c_2}.$$

The $a_1$ operator for this predicate is obviously a distributing one. In this example
$$P_1 = y^{b_1} z^{c_1} \vee y^{b_1} z^{c_2}, \text{ and } P_2 = y^{b_1} z^{c_2}, \text{ i.e. } P_1 \leftarrow P_2.$$
In order to answer the second question, it is necessary to exclude from the original equation all variables except those considered, and to investigate the resulting equation with fewer variables, describing all valid feature value sets. The work [21] considers a fairly wide class of predicates for which it is possible to specify an efficient algorithm for eliminating variables without increasing the size of the original formula. We extend here this class by adding some additional properties. Consider the following properties of the existence quantifier:

1. $\exists x\, x^a = 1$.
2. $\exists x \neg x^a = 1$.
3. $\exists x (\neg (P(x) Q(x))) = \exists x \neg P(x) \vee \exists x \neg Q(x)$.
4. $\exists x (P(x) \vee Q(x)) = \exists x P(x) \vee \exists x Q(x)$.
5. $\exists x (P(x) \& Q(y) = \exists x P(x) \& Q(y)$.
6. $\exists y (P(x) \to Q(y)) = P(x) \to \exists y Q(y)$.
7. $\exists y (P(x) \to Q(y)) = P(x) \to \exists y Q(y)$.
8. Suppose $P_i(x) \& P_j(x) = 0, i \neq j, i,j = 1,2,\ldots,k$, then:
$$\exists y \big( (P_1(x) \to Q_1(y)) \& (P_2(x) \to Q_2(y)) \& \ldots$$
$$\& (P_k(x) \to Q_k(y)) \big) = (P_1(x) \to \exists y Q_1(y)) \&$$
$$\& (P_2(x) \to \exists y Q_2(y)) \& \ldots \& (P_k(x) \to \exists y Q_k(y)).$$
9. If the identity $P_i(x) \equiv 0$ is not true for any $i = 1,2,\ldots,k$ and $P_i(x) \& P_j(x) = 0$ fo $i \neq j, i,j = 1,2,\ldots,k$, then:
$$\exists x \big( (P_1(x) \to Q_1(y)) \& (P_2(x) \to Q_2(y)) \& \ldots$$
$$\& (P_k(x) \to Q_k(y)) \big) = Q_1(y) \vee Q_2(y) \vee \ldots \vee Q_k(y).$$

The properties listed above allow describing a broad class of finite predicates (correspondingly equations) defined on the set of variables $\{x, y, \ldots, z\}$, for which it is easy to find links between selected variables without any increase in the size of the original formulas. Let us define such a class recursively.

1. All "recognitions" $x^a, x^b, \ldots, x^c$ ($a, b, \ldots, c$ — symbols belonging to the domain for the variable $x$) belong to $\Delta_x$.
2. All the negations $\neg x^a, \neg x^b, \ldots, \neg x^c$ belong to $\Delta_x$.
3. If predicates $\neg P(x), \neg Q(x)$ belong to $\Delta_x$, then the predicate $\neg(P(x)Q(x))$ belong to $\Delta_x$
4. Any predicate not depending on the variable $x$, belongs to $\Delta_x$.
5. If predicates $P_1$ and $P_2$ belong to $\Delta_x$, then the predicate $P = P_1 \vee P_2$ belongs to $\Delta_x$.
6. If the predicate $P_1$ belongs to $\Delta_x$, and the predicate $P_2$ does not depend on $x$, then the predicate $P = P_1 \& P_2$ belongs to $\Delta_x$.
7. If the predicate $P_1$ does not depend on $x$, and the predicate $P_2$ belongs to $\Delta_x$, then the predicate $P = P_1 \rightarrow P_2$ belongs to $\Delta_x$.
8. Let predicates $P_1, P_2, \ldots, P_k$ do not depend on $x$; $P_i \& P_j = 0$ for $i \neq j, i, j = 1, 2, \ldots, k$, predicates $Q_1, Q_2, \ldots, Q_k$ belong to $\Delta_x$; then
$$P = (P_1 \rightarrow Q_1) \& (P_2 \rightarrow Q_2) \& \ldots \& (P_k \rightarrow Q_k)$$
belongs to $\Delta_x$.
9. If the predicates $P_1, P_2, \ldots, P_k$ depend only on $x$, $P_i \& P_j = 0$ for $i \neq j, i, j = 1, 2, \ldots, k$; for any $i = 1, 2, \ldots, k$ the identity $P_i \equiv 0$ is not true; predicates $Q_1, Q_2, \ldots, Q_k$ do not depend on $x$; then the predicate
$$P = (P_1 \rightarrow Q_1) \& (P_2 \rightarrow Q_2) \& \ldots \& (P_k \rightarrow Q_k)$$
belongs to $\Delta_x$.

One may need also to exclude extra variables with the help of the universal quantifier. In this case we can use the following properties of this quantifier:

1. $\forall x\, x^a = 0$.
2. $\forall x \neg x^a = 0$.
3. $\forall x \neg(P(x) \vee Q(x)) = \forall x \neg P(x) \& \forall x \neg Q(x)$
4. $\forall x(P(x) \& Q(x)) = \forall x P(x) \& \forall x Q(x)$.
5. $\forall x(P(x) \vee Q(y) = \forall x P(x) \vee Q(y)$.
6. $\forall y(P(x) \& Q(y)) = P(x) \& \forall y Q(y)$.
7. Suppose

$$P_i(x) \& P_j(x) = 0, \ i \neq j, i, j = 1, 2, \ldots, k,$$

then:

$$\forall y\big((P_1(x) \& Q_1(y)) \vee (P_2(x) \& Q_2(y)) \vee \ldots$$
$$\vee (P_k(x) \& Q_k(y))\big) = (P_1(x) \& \forall y Q_1(y)) \vee$$
$$\vee (P_2(x) \& \forall y Q_2(y)) \vee \ldots \vee (P_k(x) \& \forall y Q_k(y)).$$

8. If the identity $P_i(x) \equiv 0$ is not true for any $i = 1, 2, \ldots, k$ and $P_i(x) \& P_j(x) = 0$ for $i \neq j, j = 1, 2, \ldots, k$, then:

$$\forall x\big((P_1(x) \& Q_1(y)) \vee (P_2(x) \& Q_2(y)) \vee \ldots$$
$$\vee (P_k(x) \& Q_k(y))\big) = Q_1(y) \& Q_2(y) \& \ldots \& Q_k(y).$$

We can recursively define a class of predicates $\Sigma_x$ from which it is possible to exclude the variable $x$ without an increase in the size of the formula:

1. All the "recognitions" $x^a, x^b, \ldots, x^c$ belong to $\Sigma_x$.
2. All the negations $\neg x^a, \neg x^b, \ldots, \neg x^c$ that do not depend on $x$ belong to $\Sigma_x$.
3. If $\neg P_1$ and $\neg P_2$ belong to $\Sigma_x$, then $\neg(P_1 \vee P_2)$ belongs to $\Sigma_x$.
4. If predicates $P_1$ and $P_2$ belong to $\Sigma_x$, then the predicate $P = P_1 \& P_2$ belongs to $\Sigma_x$.

5. If a $P_1$ belongs to $\Sigma_x$, and a predicate $P_2$ do not depend on $x$, then the predicate $P = P_1 \lor P_2$ belongs to $\Sigma_x$.
6. If a predicate $P_1$ does not depend on $x$, and a predicate $P_2$ belongs to $\Sigma_x$, the predicate $P = P_1 \& P_2$ belongs to $\Sigma_x$.
7. Suppose predicates $P_1, P_2, \ldots, P_k$ do not depend on $x$, $P_i \& P_j = 0$ for $i \neq j, i, j = 1, 2, \ldots, k$; predicates $Q_1, Q_2, \ldots, Q_k$ belong to $\Sigma_x$, then
$$P = (P_1 \& Q_1) \lor (P_2 \& Q_2) \lor \ldots \lor (P_k \& Q_k)$$
belongs to $\Sigma_x$.
8. If predicates $P_1, P_2, \ldots, P_k$ depend only on $x$, $P_i \& P_j = 0$ for $i \neq j, i, j = 1, 2, \ldots, k$; for any $i = 1, 2, \ldots, k$ the identity $P_i \equiv 0$ is not true, predicates $Q_1, Q_2, \ldots, Q_k$ do not depend on $x$, then the predicate $P = (P_1 \& Q_1) \lor (P_2 \& Q_2) \lor \ldots \lor (P_k \& Q_k)$ belongs to $\Sigma_x$.

## 4. Experiment and results

Let us consider a medical example and investigate links between features. The predicate variables are interconnected with systems of logic equations. Solving these equations allows attributing the objects under consideration to a certain class, which characterizes determining the risk group of a patient related to some diseases.

The plan of the experiment is as follows. We use real-world medical data and code them with the help of predicate equations. We note that although some variables can take on values "unknown", this is nevertheless a case of the closed world as "unknown" just means a value from the alphabet on which a variable is defined. Thus, every domain for any variable is closed. After we have written a system of equations with the help of experts, we start deleting variables that we consider non-salient at the moment. It does not mean that in other cased other variables will be considered as non-salient. Salient variables are those for which we want to determine logic links As an output, we obtain an equation where non-salient variables are deleted. The resulting equation is simpler than the original system, and it is possible to analyze links between salient variables in a simpler way.

If we consider the information screening of medical data for assessing the development and prevention of heart and vessel diseases [22], we can select a set of features for formalizing screening procedures. Let us consider the following features and their values:

Gender: $X_1 = \{x_1^1, x_1^2\}$, where $x_1^1$ means a woman, $x_1^1$ means a man.

Age: $X_2 = \{x_2^1, x_2^2, x_2^3\}$, where $x_2^1$ is less than 40 years, $x_2^2$ is from 40 to 50 years, $x_2^3$ is greater than 50 years.

Diabetes mellitus: $X_3 = \{x_3^1, x_3^2, x_3^3, x_3^4\}$, where $x_3^1$ – yes, $x_3^2$ – no (actual diagnosis), $x_3^3$ – no (not actual diagnosis), $x_3^4$ – unknown.

Arterial hypertension: $X_4 = \{x_4^1, x_4^2, x_4^3, x_4^4\}$, where $x_4^1$ – yes, $x_4^2$ – no (actual diagnosis), $x_4^3$ – no (not actual diagnosis), $x_4^4$ – unknown.

Kidney problems: $X_5 = \{x_5^1, x_5^2, x_5^3\}$, where $x_5^1$ – yes, $x_5^2$ – no, $x_5^3$ – unknown.

Tachycardia: $X_6 = \{x_6^1, x_6^2, x_6^3, x_6^4, x_6^5\}$, where $x_6^1$ – yes (actual diagnosis), $x_6^2$ – yes (not actual diagnosis), $x_6^3$ – no (actual diagnosis), $x_6^4$ – no (not actual diagnosis), $x_6^5$ – unknown.

Here

dity of heart and vessel diseases: $X_7 = \{x_7^1, x_7^2, x_7^3\}$, where $x_7^1$ – yes, $x_7^2$ – no, $x_7^3$ – unknown.

Smoking: $X_8 = \{x_8^1, x_8^2, x_8^3\}$, where $x_8^1$ – yes, $x_8^2$ – no, $x_8^3$ – unknown.

Alcohol problems: $X_9 = \{x_9^1, x_9^2, x_9^3\}$, where $x_9^1$ – yes, $x_9^2$ – no, $x_9^3$ – unknown.

Hypodinamia: $X_{10} = \{x_{10}^1, x_{10}^2, x_{10}^3\}$, where $x_{10}^1$ – yes, $x_{10}^2$ – no, $x_{10}^3$ – unknown.

These features allow developing a model for identifying diagnostic parameters, with the help of which it is possible to determine a group of patient health $R = \{r_1, r_2, r_3, r_4\}$, where $r_1$ is a low risk of heart and vessel diseases, $r_2$ is a moderate risk, $r_3$ is a high risk, $r_4$ is a very high risk.

For determining a health group, a set of aggregated features $Q_1 - Q_3$ can be used, where $Q_1$ is expressed in terms of $X_1$ and $X_2$, $Q_2$ is expressed in terms of $X_7$ to $X_{10}$, $Q_3$ is expressed in terms of $X_3 - X_6$.

The values of each health group and each aggregated feature is divided into four classes according to the corresponding medical technological documentation (unified clinical protocol and local protocols related to the prevention of heart and vessel diseases.

For example, for forming the feature $Q_2$, the following system of predicate equations can be formed:

$$
\begin{cases}
q_2^1 = x_7^2 x_8^2 \left( x_9^2 \vee x_9^3 (x_{10}^2 \vee x_{10}^3) \right) \vee x_7^2 x_8^3 x_9^2 x_{10}^2 \vee x_7^3 x_8^2 x_{10}^2 (x_9^2 \vee x_9^3) \\
q_2^2 = x_7^2 \left( x_8^1 (x_9^1 x_{10}^2 \vee x_9^2) \vee x_9^3 (x_8^1 x_{10}^2 \vee x_8^2 x_{10}^1) \right) \vee (x_7^2 (x_8^2 x_9^1 \vee x_8^3 x_9^2) \vee (x_7^2 x_9^3 \vee x_7^3 x_9^1) x_8^3 x_{10}^2 \vee \\
\qquad \vee (x_7^2 x_8^3 \vee x_7^3 x_8^2) x_9^1 (x_{10}^2 \vee x_{10}^3) \vee x_7^3 x_8^2 (x_9^2 \vee x_9^3) \right) (x_{10}^1 \vee x_{10}^3) \vee x_7^3 x_8^1 (x_9^2 x_{10} \vee x_9^3) \vee \\
\qquad \vee x_7^3 x_8^3 x_9^2 (x_{10}^1 \vee x_{10}^2), \\
q_2^3 = x_7^1 x_{10}^2 \left( x_8^1 x_9^2 \vee x_8^2 (x_9^1 \vee x_9^2) \right) \vee (x_7^1 x_9^3 (x_8^1 \vee x_8^2) \vee (x_7^1 x_8^3 \vee x_7^3 x_8^1) x_9^1) (x_{10}^2 \vee x_{10}^3) \vee \\
\qquad \vee x_7^1 x_8^3 (x_9^2 \vee x_9^3) \vee (x_7^2 (x_8^1 (x_9^1 \vee x_9^3) \vee x_8^3 x_9^3) \vee (x_7^2 x_8^3 \vee x_7^3 x_8^2) x_9^1 x_{10}^1 \vee \\
\qquad \vee x_7^3 (x_8^1 x_9^2 \vee x_8^3 x_9^1)) (x_{10}^1 \vee x_{10}^3) \vee x_7^3 x_8^3 (x_9^2 x_{10}^3 \vee x_9^3), \\
q_2^4 = x_7^1 x_9^3 x_{10}^1 (x_8^1 \vee x_8^2) \vee (x_7^1 x_8^3 \vee x_7^3 x_8^1) x_9^1 x_{10}^1 \vee \left( x_7^1 x_8^2 x_9^1 \vee x_7^1 x_9^2 (x_8^1 \vee x_8^2) \right) (x_{10}^1 \vee x_{10}^3) \vee x_7^1 x_8^1 x_9^1.
\end{cases}
$$

The final classification can be expressed by the following system:

$$
\begin{cases}
r_1 = q_1^1 q_2^1 (q_3^1 \vee q_3^2) \vee (q_1^1 q_2^2 \vee (q_1^2 \vee q_1^3) q_2^1) q_3^1, \\
r_2 = q_1^1 (q_2^1 q_3^3 \vee q_2^2 q_3^2) \vee \left( q_1^1 (q_2^3 \vee q_2^4) \vee q_1^2 (q_2^2 \vee q_2^3) \vee q_1^3 q_2^2 \vee q_1^4 (q_2^1 \vee q_2^2) \right) (q_3^1 \vee q_3^2) \vee \\
\qquad \vee (q_1^2 \vee q_1^3) q_2^1 (q_3^2 \vee q_3^3) \vee (q_1^2 q_2^4 \vee (q_1^3 \vee q_1^4) q_2^3) q_3^1, \\
r_3 = q_2^1 q_3^4 \vee (q_1^1 \vee q_1^2 \vee q_1^3) (q_2^2 \vee q_2^3) (q_3^3 \vee q_3^4) \vee q_1^3 q_2^3 (q_2^3 \vee q_2^4) \vee (q_1^3 \vee q_1^4) q_2^4 q_3^1 \vee \\
\qquad \vee \left( q_1^1 q_2^4 \vee q_1^4 (q_2^1 \vee q_2^2) \right) q_3^3 \vee (q_1^2 q_2^4 \vee q_1^4 q_2^3) (q_3^2 \vee q_3^3), \\
r_4 = (q_1^1 \vee q_1^2) q_2^4 q_3^4 \vee q_1^3 q_2^4 (q_3^3 \vee q_3^4) \vee q_1^4 q_3^3 (q_2^2 \vee q_2^3) \vee q_1^4 q_2^2 (q_3^2 \vee q_3^3 \vee q_3^4)
\end{cases}
$$

Let us investigate logic links between discrete features $x_1 - x_{10}$. First of all, let us rewrite the system of predicate equations in the following form (2):

$$
P(q_2, x_1, \dots, x_{10}) = q_2^1 (x_7^2 x_8^2 \left( x_9^2 \vee x_9^3 (x_{10}^2 \vee x_{10}^3) \right) \vee x_7^2 x_8^3 x_9^2 x_{10}^2 \vee x_7^3 x_8^2 x_{10}^2 (x_9^2 \vee x_9^3)) \vee
$$
$$
\vee q_2^2 (x_7^2 \left( x_8^1 (x_9^1 x_{10}^2 \vee x_9^2) \vee x_9^3 (x_8^1 x_{10}^2 \vee x_8^2 x_{10}^1) \right) \vee (x_7^2 (x_8^2 x_9^1 \vee x_8^3 x_9^2) \vee (x_7^2 x_9^3 \vee x_7^3 x_9^1) x_8^3 x_{10}^2 \vee
$$
$$
\vee (x_7^2 x_8^3 \vee x_7^3 x_8^2) x_9^1 (x_{10}^2 \vee x_{10}^3) \vee x_7^3 x_8^2 (x_9^2 \vee x_9^3)) (x_{10}^1 \vee x_{10}^3) \vee x_7^3 x_8^1 (x_9^2 x_{10}^2 \vee x_9^3) \vee
$$
$$
\vee x_7^3 x_8^3 x_9^2 (x_{10}^1 \vee x_{10}^2)) \vee
$$
$$
\vee q_2^3 (x_7^1 x_{10}^2 \left( x_8^1 x_9^2 \vee x_8^2 (x_9^1 \vee x_9^2) \right) \vee (x_7^1 x_9^3 (x_8^1 \vee x_8^2) \vee (x_7^1 x_8^3 \vee x_7^3 x_8^1) x_9^1) (x_{10}^2 \vee x_{10}^3) \vee
$$
$$
\vee x_7^1 x_8^3 (x_9^2 \vee x_9^3) \vee (x_7^2 (x_8^1 (x_9^1 \vee x_9^3) \vee x_8^3 x_9^3) \vee (x_7^2 x_8^3 \vee x_7^3 x_8^2) x_9^1 x_{10}^1 \vee
$$
$$
\vee x_7^3 (x_8^1 x_9^2 \vee x_8^3 x_9^1)) (x_{10}^1 \vee x_{10}^3) \vee x_7^3 x_8^3 (x_9^2 x_{10}^3 \vee x_9^3)) \vee
$$
$$
\vee q_2^4 (x_7^1 x_9^3 x_{10}^1 (x_8^1 \vee x_8^2) \vee (x_7^1 x_8^3 \vee x_7^3 x_8^1) x_9^1 x_{10}^1 \vee \left( x_7^1 x_8^2 x_9^1 \vee x_7^1 x_9^2 (x_8^1 \vee x_8^2) \right) (x_{10}^1 \vee x_{10}^3) \vee
$$
$$
x_7^1 x_8^1 x_9^1) = 1.
$$

It can be seen that this predicate belongs to the class $\Delta_{x7}$. Let us investigate the link between all variables except $x_7$. This elimination will give us the link between the variables $q_2, x_1, \dots, x_6, x_8, x_9, x_{10}$:

$$
F = \exists x_7 P(q_2, x_1, \dots, x_{10}) =
$$
$$
= q_2^1 (x_8^2 \left( x_9^2 \vee x_9^3 (x_{10}^2 \vee x_{10}^3) \right) \vee x_8^3 x_9^2 x_{10}^2 \vee x_8^2 x_{10}^2 (x_9^2 \vee x_9^3)) \vee
$$

$$\vee\, q_2^2 \left( x_8^1(x_9^1 x_{10}^2 \vee x_9^2) \vee x_9^3(x_8^1 x_{10}^2 \vee x_8^2 x_{10}^1) \right) \vee \left( (x_8^2 x_9^1 \vee x_8^3 x_9^2) \vee (x_9^3 \vee x_9^1)x_8^3 x_{10}^2 \vee \right.$$
$$\vee\, (x_8^3 \vee x_8^2)x_9^1(x_{10}^2 \vee x_{10}^3) \vee x_8^2(x_9^2 \vee x_9^3) \big)(x_{10}^1 \vee x_{10}^3) \vee x_8^1(x_9^2 x_{10}^2 \vee x_9^3) \vee$$
$$\vee\, x_8^3 x_9^2(x_{10}^1 \vee x_{10}^2)) \vee$$
$$\vee\, q_2^3(x_{10}^2 \left( x_8^1 x_9^2 \vee x_8^2(x_9^1 \vee x_9^2) \right) \vee (x_9^3(x_8^1 \vee x_8^2) \vee (x_8^3 \vee x_7^3 x_8^1)x_9^1)(x_{10}^2 \vee x_{10}^3) \vee$$
$$\vee\, x_8^3(x_9^2 \vee x_9^3) \vee \left( (x_8^1(x_9^1 \vee x_9^3) \vee x_8^3 x_9^3) \vee (x_8^3 \vee x_8^2)x_9^1 x_{10}^1 \vee \right.$$
$$\vee\, (x_8^1 x_9^2 \vee x_8^3 x_9^1) \big)(x_{10}^1 \vee x_{10}^3) \vee x_8^3(x_9^2 x_{10}^3 \vee x_9^3)) \vee$$
$$\vee\, q_2^4(x_9^3 x_{10}^1(x_8^1 \vee x_8^3) \vee (x_8^3 \vee x_8^1)x_9^1 x_{10}^1 \vee \left( x_8^2 x_9^1 \vee x_9^2(x_8^1 \vee x_8^2) \right)(x_{10}^1 \vee x_{10}^3) \vee x_8^1 x_9^1) = 1.$$

It should be noted that the size of the original formula has not increased, which is because the predicate $P(q_2, x_1, \dots, x_{10})$ belongs to $\Delta_{x7}$ .

Suppose we are interested in the link between $q_2, x_9, x_{10}$. Let us eliminate the other features from the predicate $F(x_1, \dots, x_{10})$:

$$G(q_2, x_9, x_{10}) = \exists x_1 \exists x_2 \exists x_3 \exists x_4 \exists x_5 \exists x_6 \exists x_7 \exists x_8 P(q_2, x_1, \dots, x_{10}) =$$
$$= q_2^1(\left( x_9^2 \vee x_9^3(x_{10}^2 \vee x_{10}^3) \right) \vee x_9^2 x_{10}^2 \vee x_{10}^2(x_9^2 \vee x_9^3)) \vee$$
$$\vee\, q_2^2 \left( (x_9^1 x_{10}^2 \vee x_9^2) \vee x_9^3(x_{10}^2 \vee x_{10}^1) \right) \vee \left( (x_9^1 \vee x_9^2) \vee (x_9^3 \vee x_9^1)x_{10}^2 \vee \right.$$
$$\vee\, x_9^1(x_{10}^2 \vee x_{10}^3) \vee (x_9^2 \vee x_9^3) \big)(x_{10}^1 \vee x_{10}^3) \vee (x_9^2 x_{10}^2 \vee x_9^3) \vee$$
$$\vee\, x_9^2(x_{10}^1 \vee x_{10}^2)) \vee$$
$$\vee\, q_2^3(x_{10}^2(x_9^1 \vee x_9^2) \vee x_9^3(x_{10}^2 \vee x_{10}^3) \vee$$
$$\vee\, (x_9^2 \vee x_9^3) \vee \left( (x_9^1 \vee x_9^3) \vee x_9^1 x_{10}^1 \vee \right.$$
$$\vee\, (x_9^2 \vee x_9^1) \big)(x_{10}^1 \vee x_{10}^3) \vee (x_9^2 x_{10}^3 \vee x_9^3)) \vee$$
$$\vee\, q_2^4(x_9^3 x_{10}^1 \vee x_9^1 x_{10}^1 \vee (x_9^1 \vee x_9^2)(x_{10}^1 \vee x_{10}^3) \vee x_9^1) = 1.$$

Again, we have reduced the original formula and obtained a simpler dependence between selected medical features. After the necessary dependence is obtained, we can solve the resulting equation with one or several target variables.

## 5. Conclusions and further research

In this paper finite predicate equations of different types have been considered. The description of classification problems based on predicate equations have been presented. The problem of object classification on the basis of features taking on discrete values has been described mathematically as a solution of predicate equations. A broad class of predicates from which it is possible to delete extra variables and focus on links between salient variables has been described. A method for deleting non-salient variables by the application of the existential quantifier has been suggested and demonstrated on a real-world medical example.

Although some variables in the medical example can take on values "unknown", we deal with the closed world here as "unknown" means just an element from the domain for the variable. All the domains are strictly defined and cannot be completed with any other elements. The main advantage of this method based on a specific structure of predicate systems lies in the fact that after deleting non-salient variables the original system (or equation) is simplified, which is due to special properties of quantifiers. Salient variables are not necessarily fixed forever. It is up to researcher to decide what logic connections are important at the moment.

As a further research direction, we are going to extend classes of predicates from which it is easy to delete non-salient variables Such classes are much more complicated than relational structures, but many practical problems require a corresponding knowledge representation and analysis.

## 6. References

[1] G. A. Oparin, V. G. Bogdanova and A. A. Pashinin, "Classification in Binary Feature Space Using Logical Dynamic Models," 44th International Convention on Information, Communication and Electronic Technology (MIPRO), Opatija, Croatia, 2021, pp. 1020-1025, doi: 10.23919/MIPRO52101.2021.9596697.

[2] A. Kabulov, E. Urunboev and I. Saymanov, "Object recognition method based on logical correcting functions," 2020 International Conference on Information Science and Communications Technologies (ICISCT), Tashkent, Uzbekistan, 2020, pp. 1-4, doi: 10.1109/ICISCT50599.2020.9351473.

[3] A. M. Ahmed, S. K. Ibrahim and S. Yacout, "Hyperspectral Image Classification Based on Logical Analysis of Data," 2019 IEEE Aerospace Conference, Big Sky, MT, USA, 2019, pp. 1-9, doi: 10.1109/AERO.2019.8742023.

[4] S. Kim, S. Noh and H. S. Ryoo, "Identifying Combinatorial Significance for Classification of Alzheimer's Disease Proteomics Expression with Logical Analysis of Data," 2021 IEEE International Conference on Bioinformatics and Biomedicine (BIBM), Houston, TX, USA, 2021, pp. 1661-1663, doi: 10.1109/BIBM52615.2021.9669835.

[5] I. Povkhan and M. Lupei, "The Algorithmic Classification Trees," 2020 IEEE Third International Conference on Data Stream Mining & Processing (DSMP), Lviv, Ukraine, 2020, pp. 37-43, doi: 10.1109/DSMP47368.2020.9204198.

[6] M. Tamilselvi, G. Ramkumar, G. Anitha, P. Nirmala and S. Ramesh, "A Novel Text Recognition Scheme using Classification Assisted Digital Image Processing Strategy," 2022 International Conference on Advances in Computing, Communication and Applied Informatics (ACCAI), Chennai, India, 2022, pp. 1-6, doi: 10.1109/ACCAI53970.2022.9752542.

[7] S. N. Wagle and B. Kovalerchuk, "Interactive Visual Self-service Data Classification Approach to Democratize Machine Learning," 2020 24th International Conference Information Visualisation (IV), Melbourne, Australia, 2020, pp. 280-285, doi: 10.1109/IV51561.2020.00052.

[8] B. Jean-Marc and O. Lafitte, "Combining weak classifiers: a logical analysis," 2021 23rd International Symposium on Symbolic and Numeric Algorithms for Scientific Computing (SYNASC), Timisoara, Romania, 2021, pp. 178-181, doi: 10.1109/SYNASC54541.2021.00038.

[9] Z. Chen, A. Xu, Y. Zhou and Y. Gai, "Research on Pulsar Classification Based on Machine Learning," 2020 3rd International Conference on Intelligent Autonomous Systems (ICoIAS), Singapore, 2020, pp. 14-18, doi: 10.1109/ICoIAS49312.2020.9081836.

[10] C. Nwankwo, H. Wimmer, L. Chen and J. Kim, "Text Classification of Digital Forensic Data," 2020 11th IEEE Annual Information Technology, Electronics and Mobile Communication Conference (IEMCON), Vancouver, BC, Canada, 2020, pp. 0661-0667, doi: 10.1109/IEMCON51383.2020.9284913.

[11] I. Shubin, S. Snisar and S. Litvin, "Categorical Analysis of Logical Networks in Application to Intelligent Radar Systems," 2020 IEEE International Conference on Problems of Infocommunications. Science and Technology (PIC S&T), Kharkiv, Ukraine, 2020, pp. 235-238, doi: 10.1109/PICST51311.2020.9467893.

[12] D. Cenzer, V. W. Marek and J. B. Remmel, "On the complexity of index sets for finite predicate logic programs which allow function symbols," in Journal of Logic and Computation, vol. 30, no. 1, pp. 107-156, Jan. 2020, doi: 10.1093/logcom/exaa005.

[13] A. Kabulov, E. Urunbayev and A. Ashurov, "Logic Method of Finding Maximum Joint Subsystems of Systems of Boolean Equations," 2020 International Conference on Information Science and Communications Technologies (ICISCT), Tashkent, Uzbekistan, 2020, pp. 1-5, doi: 10.1109/ICISCT50599.2020.9351394.

[14] H. Qi, B. Li, R. -J. Jing, A. Proutiere and G. Shi, "Distributedly Solving Boolean Equations over Networks," 2020 59th IEEE Conference on Decision and Control (CDC), Jeju, Korea (South), 2020, pp. 560-565, doi: 10.1109/CDC42340.2020.9304144.

[15] F. N. Castro, O. E. González and L. A. Medina, "Diophantine Equations With Binomial Coefficients and Perturbations of Symmetric Boolean Functions," in IEEE Transactions on Information Theory, vol. 64, no. 2, pp. 1347-1360, Feb. 2018, doi: 10.1109/TIT.2017.2750674.

[16] L. Gong, J. Zhang, L. Sang, H. Liu and Y. Wang, "The Unpredictability Analysis of Boolean Chaos," in IEEE Transactions on Circuits and Systems II: Express Briefs, vol. 67, no. 10, pp. 1854-1858, Oct. 2020, doi: 10.1109/TCSII.2019.2949571.

[17] T. Kosovskaya, "Implementation of Formula Partial Sequence for Rough Solution of AI Problems in the Framework of the Logic-Predicate Approach," 2019 Computer Science and Information Technologies (CSIT), Yerevan, Armenia, 2019, pp. 65-68, doi: 10.1109/CSITechnol.2019.8895153.

[18] Y. -F. Tseng and S. -J. Gao, "Efficient Subset Predicate Encryption for Internet of Things," 2021 IEEE Conference on Dependable and Secure Computing (DSC), Aizuwakamatsu, Fukushima, Japan, 2021, pp. 1-2, doi: 10.1109/DSC49826.2021.9346245.

[19] K. Smelyakov, A. Chupryna, O. Bohomolov and N. Hunko, "The Neural Network Models Effectiveness for Face Detection and Face Recognition," 2021 IEEE Open Conference of Electrical, Electronic and Information Sciences (eStream), 2021, pp. 1-7, doi: 10.1109/eStream53087.2021.9431476.

[20] Sharonova N. Issues of Fact-based Information Analysis [Electronic resource] / N. Sharonova, A. Doroshenko, O. Cherednichenko // Computational linguistics and intelligent systems (COLINS 2018) : proc. of the 2nd Intern. Conf., June 25-27, 2018. Vol. 1: Main Conference / ed.: V. Lytvyn [et al.]. – Electron. text data. – Lviv, 2018. – P. 11-19. – URL: http://ceur-ws.org/Vol-2136/10000011.pdf, (accessed 01.06.2020).

[21] D.E.Sitnikov, B D'Cruz, P.E. Sitnikova, "Discovering salient data features by composing and manipulating logical equations", Data Mining II, WIT Press, 2000, 241-248.

[22] Melnik K. Towards medical screening information technology: the healthgrid-based approach / K. Melnik, O. Cherednichenko, V. Glushko // Information Systems: Methods, Models, and Applications. – Heidelberg : Springer, 2013. – P. 202-204.