



# Interactive error correction in implicative theories



Sergei O. Kuznetsov<sup>a</sup>, Artem Revenko<sup>a,b,\*</sup>

<sup>a</sup> National Research University Higher School of Economics, Pokrovskiy bd. 11, 109028 Moscow, Russia

<sup>b</sup> Technische Universität Dresden, Zellescher Weg 12-14, 01069 Dresden, Germany

## ARTICLE INFO

### Article history:

Received 20 October 2014

Received in revised form 3 April 2015

Accepted 10 June 2015

Available online 16 June 2015

### Keywords:

Implicative theory

Error correction

Closure system

Formal concept analysis

## ABSTRACT

Errors in implicative theories coming from binary data are studied. First, two classes of errors that may affect implicative theories are singled out. Two approaches for finding errors of these classes are proposed, both of them based on methods of Formal Concept Analysis. The first approach uses the cardinality minimal (canonical or Duquenne–Guigues) implication base. The construction of such a base is computationally intractable. Using an alternative approach one checks possible errors on the fly in polynomial time via computing closures of subsets of attributes. Both approaches are interactive, based on questions about the validity of certain implications. Results of computer experiments are presented and discussed.

© 2015 Elsevier Inc. All rights reserved.

## 1. Introduction

### 1.1. Motivation

Implicative theories consisting of formulas of the form “if  $A$ , then  $B$ ” provide a standard way for describing the structure of domain knowledge. They are extensively used in various research areas, e.g., biology [18], pharmacology [6,5], semantic web [19], knowledge discovery [12,34], decision making [26], classification [24], ontology engineering [2]. In many cases the exactness of rules plays a crucial role, for example in research related to strictly formalized domains like Boolean algebras [22], algebraic lattices [7], or algebraic identities [27].

In many applications an exact implicative theory is constructed from a piece of available data. It is well known that a single mistake in this data can drastically change the resulting implicative theory [14] (the same is true for association rules if there are some exceptions and an error). The implicative theory is not going to recover from this error even if further error-free data is added to the underlying set. Therefore, implicative theories are not error tolerant. However, in the real-world applications, especially if multiple users are expected to work with data, one cannot guarantee the absence of errors. More than that, someone may be willing to spoil the result on purpose by adding erroneous instances, in order to prevent from discovering valid implications. Therefore, a procedure for recovering from errors is essential for the usage of implicative theories.

Here we assume that in the beginning there is already some data on hands and new data arrives in the work flow. The goal is to guarantee the correctness of the implicative theory with respect to the initial data which are considered to be reliable. We do not assume that a user, which is going to work with the data and the implicative theory, is always able to explicitly state any knowledge about data domain or has any knowledge about methods in use. That is why it is

\* Corresponding author at: Technische Universität Dresden, Zellescher Weg 12-14, 01069 Dresden, Germany.

E-mail addresses: skuznetsov@hse.ru (S.O. Kuznetsov), artem\_viktorovich.revenko@mailbox.tu-dresden.de (A. Revenko).

Name	Sex	Year of Birth	Lawful Age	...
Helga	f	1995	n	...
Daria	f	1980	y	...
Patrick	m	1986	y	...
John	m	1996	n	...
↑				
George	m	1980	n	...

Fig. 1. Data table and new entry from Example 1.

important to develop a transparent and easy method for error correction. In particular, it is important to find and output possible errors in a human understandable form. To attain this goal a natural framework can be that of Formal Concept Analysis (FCA) [14], where methods and algorithms for finding implicative theories of binary data (formal contexts) are well elaborated and widely used [13,30].

**Example 1.** To illustrate our ideas we provide a use-case example. Let there be data from Fig. 1 on hands. New data is coming from an untrusted source and it is intended to be added to the existing data. The user expects possible errors in new data, however, the user is not able to check every single entry (possibly, due to a large number of columns). The solution we propose in this paper would output the question: Does ‘Year of Birth: 1980’ imply ‘Lawful Age’? As we are now in year 2015, the answer is obviously ‘Yes’ and, therefore, an error is revealed.

## 1.2. Related work

Methods for imputing missing values are well studied. In [33] and [31] detailed overviews of existing techniques are presented. Among others there are techniques of ignoring entries with missing values, imputing average values, and more complicated ones such as decision trees, neural networks [31], Nearest Neighbor approach [16]. Having a missing value, there is no need to search for an error, as it is clear from the problem statement which value should be changed (or imputed). An approach proposed in this paper bares some similarity to the Nearest Neighbor method, but aims at solving a different task. Besides that, the imputation techniques (like, e.g. averaging) are mostly not relevant for binary data.

Error finding and eliminating are widely discussed in various fields of computer science. The problems of lineage or data provenance, where one needs to explain errors, trace reasons for a query, etc. are well known in KDD domain [32]. These techniques are very useful and efficient, however, they are not appropriate for correcting errors in binary data tables.

In [9] an impressive way of using expert knowledge presented in the form of editing rules and certain regions for databases are surveyed. Information in the form of editing rules prevents the errors from getting in to the database. The approach presented in this paper aims at finding and correcting errors without any previously formalized knowledge.

The paper [10] presents an interesting approach to dealing with mistakes in answering questions (like the ones we will discuss below) in the process of knowledge base completion within the framework of Description Logics. This approach allows recovering from such mistakes in such an effective manner that the information input is used upon mistake recovery. However, the detection and correction of mistakes is left to pinpointing.

Pinpointing is a very helpful technique for recovering from inconsistencies. The goal of pinpointing is the following: for a given inconsistent set of rules (not only implicative) find minimal inconsistent subsets [3,23]. The inconsistency is detected via checking if a certain erroneous consequence holds. This technique is successfully applied in different description logics. The complexity of pinpointing is normally beyond polynomial. An approach introduced in this paper (Section 4, base approach) is closely related to pinpointing; it proceeds from knowledge base constructed from data. The complexity is also beyond polynomial. However, an alternative approach (Section 4, closure approach) takes the advantage of having the data and proposes a polynomial-time solution. In this work we do not modify the knowledge base directly, but we correct the errors in data in such a way that the corresponding implicative theory becomes error-free.

As implicative theories is another view of Horn theories [14], the problem of finding explanations in Horn theories turns out to be closely connected to our problem. Namely, an entry in the binary data table can also be considered as a fact to be explained. In [17] it is shown that such explanations may be found in polynomial time. However, here we aim at explaining existence or absence of all attributes at the same time. Also we state our task and our solutions in a different language and provide algorithms for practical usage. The case of negative attributes is not covered in [17] as opposed to this work.

The present paper is a follow-up work to [28].

**Remark 1.** In this paper we assume that we can put questions to an expert in the domain who gives correct answers.

**Remark 2.** All sets and contexts we consider in this paper are assumed to be finite, which practically means an obvious constraint on finiteness of data at hand.

### 1.3. Contributions

We introduce an interactive procedure for making implicative theory of data error-free. This goal is achieved via finding and eliminating errors in rows of data tables. In FCA terms, we propose an approach for finding errors in descriptions of new objects (intents in terms of FCA) that affect the canonical implication base.

1. We introduce two possible classes of errors in data tables (formal contexts, Section 3).
2. We introduce two approaches to finding and eliminating errors of certain classes (Sections 4). Both aim at restoring dependencies from data domain and eliminating errors in implicative theory.
  - (a) One approach is based on finding those implications from an implication base that are not respected by the new object intent (base approach). However, the base approach leads to an intractable solution, because constructing an implication base is intractable.
  - (b) We introduce another approach (closure approach), where we do not need to compute the set of all implications, and prove its effectiveness. We show that it helps to find all possible errors of certain types (Proposition 1) in polynomial time (Proposition 2).

The approaches are experimentally compared in Section 5.

## 2. Main definitions

In what follows we keep to standard definitions of FCA [14]. Let  $G$  and  $M$  be sets and let  $I \subseteq G \times M$  be a binary relation between  $G$  and  $M$ . The triple  $\mathbb{K} := (G, M, I)$  is called a (*formal*) *context*. Set  $G$  is called a set of *objects*, set  $M$  is called a set of *attributes*,  $I$  is called *incidence relation*.

Consider mappings  $\varphi: 2^G \rightarrow 2^M$  and  $\psi: 2^M \rightarrow 2^G$ :  $\varphi(X) := \{m \in M \mid gIm \text{ for all } g \in X\}$ ,  $\psi(A) := \{g \in G \mid gIm \text{ for all } m \in A\}$ . Mappings  $\varphi$  and  $\psi$  define a *Galois connection* between  $(2^G, \subseteq)$  and  $(2^M, \subseteq)$ , i.e.  $\varphi(X) \subseteq A \Leftrightarrow \psi(A) \subseteq X$ . Hence, for any  $X_1, X_2 \subseteq G$ ,  $A_1, A_2 \subseteq M$  one has

1.  $X_1 \subseteq X_2 \Rightarrow \varphi(X_2) \subseteq \varphi(X_1)$
2.  $A_1 \subseteq A_2 \Rightarrow \psi(A_2) \subseteq \psi(A_1)$
3.  $X_1 \subseteq \psi\varphi(X_1)$  and  $A_1 \subseteq \varphi\psi(A_1)$

Usually, instead of  $\varphi$  and  $\psi$  a single notation  $(\cdot)'$  is used.  $(\cdot)'$  is usually called a *derivation operator*. For  $X \subseteq G$  the set  $X'$  is called the *intent* of  $X$ . Similarly, for  $A \subseteq M$  the set  $A'$  is called the *extent* of  $A$ . Operator  $(\cdot)''$  is idempotent, extensive and monotone, i.e., has the properties of algebraical closure both on  $2^G$  and  $2^M$ . Hence,  $(Z)''$  is called *closure* of  $Z$  in  $\mathbb{K}$  for  $Z \subseteq M$  or  $Z \subseteq G$ . If  $(Z)'' = Z$ , set  $Z$  is called *closed* in  $\mathbb{K}$ . Applying Properties 1 and 2 consequently one gets the *monotonicity* property: for any  $Z_1, Z_2 \subseteq G$  or  $Z_1, Z_2 \subseteq M$  one has  $Z_1 \subseteq Z_2 \Rightarrow Z_1'' \subseteq Z_2''$ .

In [29] authors introduce a generalized framework for considering positive and negative attributes. In this paper we also introduce negative attributes, however, we do not need the whole framework for our purpose. Our definitions comply with the definitions from [29].

The set  $\bar{M} := \{\bar{m} \mid m \in M\}$  is called the set of *negative attributes*. Consider the following relation  $\bar{I} := \{(g, \bar{m}) \mid (g, m) \in (G \times M) \setminus I\}$  between  $G$  and  $M$ . The context  $\mathbb{K}^\delta := (G, M \cup \bar{M}, I \cup \bar{I})$  is called the *dichotomized context* to  $\mathbb{K}$ , the corresponding derivation operator is denoted by  $(\cdot)^\delta$ . Let  $X \subseteq G$ . Note that  $\bar{m} \in X^\delta$  iff  $m \notin g'$  for all  $g \in X$ . If  $\bar{m} \in X^\delta$  then, as it does not lead to ambiguity, we informally write  $\bar{m} \in X'$ . In this paper objects and context are represented without negative attributes, however, in the processing stage they are normally converted to the dichotomized representation in order to be able to work with negative attributes.

Consider the context  $\mathbb{K}^\delta = (G, \bar{M} \cup \bar{M}, \bar{I} \cup \bar{I})$ . This context is isomorphic to the context  $\mathbb{K}^\delta = (G, M \cup \bar{M}, I \cup \bar{I})$  and  $\bar{m} \in X^\delta \Leftrightarrow m \in X^\delta$ .

A *formal concept* of a formal context  $(G, M, I)$  is a pair  $(X, A)$ , where  $X \subseteq G$ ,  $A \subseteq M$ ,  $X' = A$ , and  $A' = X$ . The set  $X$  is called the *extent*, and the set  $A$  is called the *intent* of the concept  $(X, A)$ .

One says that an object  $g$  such that  $g' \neq \emptyset$  is *reducible* in a context  $(G, M, I)$  iff  $\exists X \subseteq G \setminus \{g\} : g' = \bigcap_{j \in X} j'$ . Removing

reducible objects does not change the concept lattice up to isomorphism.

In this paper implicative theories are formalized in terms of implication bases. An *implication* of  $\mathbb{K} := (G, M, I)$  is defined as a pair  $(A, B)$ , written  $A \rightarrow B$ , where  $A, B \subseteq M$ .  $A$  is called the *premise*,  $B$  is called the *conclusion* of implication  $A \rightarrow B$ . Implication  $A \rightarrow B$  is *respected by a set of attributes*  $N$  if  $A \not\subseteq N$  or  $B \subseteq N$ . Implication  $A \rightarrow B$  holds (is valid) in  $\mathbb{K}$  if it is respected by all  $g', g \in G$ , i.e. every object, that has all the attributes from  $A$ , also has all the attributes from  $B$ , or, equivalently, if  $A' \subseteq B'$ . Implications satisfy *Armstrong rules*:

$$\frac{}{A \rightarrow A} \quad , \quad \frac{A \rightarrow B}{A \cup C \rightarrow B} \quad , \quad \frac{A \rightarrow B, B \cup C \rightarrow D}{A \cup C \rightarrow D}$$

*Support* of implication  $A \rightarrow B$  in context  $\mathbb{K}$  is  $(A \cup B)'$ , i.e., the set of all objects of  $\mathbb{K}$ , whose intents contain the premise and the conclusion of the implication. A *unit implication* is defined as an implication with only one attribute in conclusion,

i.e.  $A \rightarrow b$ , where  $A \subseteq M$ ,  $b \in M$ . Using Armstrong rules, every implication  $A \rightarrow B$  can be represented as a set of *unit implications*  $\{A \rightarrow b \mid b \in B\}$ , so one can always observe only unit implications without loss of generality.

Consider implications of the form  $A \rightarrow \bar{b}$ , where  $A \subseteq M$ ,  $\bar{b} \in \bar{M}$  in the dichotomized context  $\mathbb{K}^\delta$ . This implication is said to be respected by  $N \subseteq M$  if  $A \not\subseteq N$  or  $b \in M \setminus N$ . This implication holds in  $\mathbb{K}^\delta$  iff  $A^\delta \subseteq \bar{b}^\delta$ . In this paper all the implications with negative attributes are considered as implications of the dichotomized context.

An *implication base* of a context  $\mathbb{K}$  is defined as a set  $\mathcal{I}$  of implications of  $\mathbb{K}$ , from which any valid implication for  $\mathbb{K}$  can be deduced by Armstrong rules and none of the proper subsets of  $\mathcal{I}$  has this property. A cardinality minimal implication base was characterized in [15] and is known as the *canonical implication base*, or *Duquenne–Guigues base*, or *stembase*. In [13] the premises of implications of the canonical base were characterized in terms of pseudo-intents. A subset of attributes  $P \subseteq M$  is called a *pseudo-intent* if  $P \neq P''$  and for every pseudo-intent  $Q$  such that  $Q \subset P$ , one has  $Q'' \subset P$ . The canonical implication base looks as follows:  $\{P \rightarrow (P'' \setminus P) \mid P - \text{pseudo-intent}\}$ .

### 3. Errors in implicative theories

Without loss of generality we consider all observable properties to be expressed in terms of positive attributes from  $M$ . We aim at restoring valid implications and, therefore, correct errors in implicative theory of data. The goal is achieved if all implications are valid implications of the context. As already mentioned, all implications are reduced to unit ones.

Consider the following possible classes of implicative formulas ( $A \subseteq M$ ,  $b \in M$ ), which will be called dependencies:

1. If an entity has all attributes from  $A$ , then it has attribute  $b$  ( $A \rightarrow b$ );
2. If an entity has all attributes from  $A$ , then it does not have attribute  $b$  ( $A \rightarrow \bar{b}$ );

**Remark 3.** In this work we consider only data domain dependencies in the form of implications with no negative attributes in the premise. It is possible to consider negative attributes in the premise by means of considering complementary context  $(G, M, (G \times M) \setminus I)$ . However, this is equivalent to introducing disjunction to our language:  $A \rightarrow B \vee C \Leftrightarrow A, \bar{B} \rightarrow C$ . Then, having negation and disjunction we end up in the full propositional logic, for which computing the closure is not polynomial anymore. Therefore, it would not be possible to introduce a polynomial solution of this problem.

Only formulas of Class 1 are standard FCA implications, formulas of Class 2 are FCA implications if the negation of attributes are explicitly introduced in the context. If there are no errors in a context, all the dependencies of Class 1 are deducible from an implication base. However, if not enough data is added to the context yet, we may get false consequences. Therefore, not all valid implications of the context have to necessarily be data domain dependencies. Nevertheless, it is guaranteed that none of valid dependencies is lost, and, as new objects are added, the number of false consequences is reduced (this is essentially the idea behind Attribute Exploration [14]). The situation is different if an erroneous object (data table row) is added. The erroneous object may violate a data domain dependency. In this case, until the error is found and corrected, we are not able to deduce all dependencies valid in the data domain from the implication base, no matter how many error-free objects are added afterwards.

### 4. Finding errors

We introduce two different approaches to finding errors. The first one is based on inspecting the canonical base of a context (base approach). When adding a new object to the context one may find all implications from the canonical base of the context such that the implications are not respected by the intent of the new object. These implications are then output as questions to an expert in form of implications. If at least one of these implications is accepted, the object intent is erroneous. Since the canonical base is the most compact (in the number of implications) representation of all valid implications of a context, it is guaranteed that minimal number of questions is asked and no valid dependencies of Class 1 are left out.

This approach is related to the procedure of pinpointing. Namely, in pinpointing an invalid implication is given and the task is to find minimal subsets of rules such that the invalid implication is deducible (if any rule is removed from a minimal subset then the invalid implication is not deducible anymore). In our setting we can consider an object intent as a conjunction of attributes. The task is to find the maximal subsets of implications that are respected by the object intent, or finding the minimal subsets of implications that are violated. Therefore, in contrast to pinpointing, the statement (the conjunction of attributes) is not a consequence of the minimal subsets of implication. On the contrary, the statement violates every subset. Afterwards we take the union of all minimal violated subsets and obtain a subset of implications where each implication is violated by the new object.

Although this approach allows one to reveal all dependencies of Class 1, there are several issues. The problem of producing the canonical basis with known algorithms is intractable. Recent theoretical results [20,8,21,4] suggest that the canonical base can hardly be computed with better worst-case complexity than that of the existing approaches [13]. One can use other bases (for example, there has been recent progress in computing proper premises [30]), but the algorithms known so far are still too costly and non-minimal bases do not guarantee that the expert is asked minimal sufficient number of questions.

However, since we are only interested in implications corresponding to one object at a time, it may be not necessary to compute the whole implication base. The second approach takes this fact into account. Let  $A \subseteq M$  be the intent of the object under inspection; we separate it from the context.  $m \in A''$  iff  $\forall g \in G : A \subseteq g' \Rightarrow m \in g'$ , in other words,  $A''$  contains the attributes common to all object intents containing  $A$ . The set of unit implications  $\{A \rightarrow b \mid b \in A'' \setminus A\}$  can then be shown to the expert. If all implications are rejected, no attributes are forgotten in the new object intent. Otherwise, there are missing attributes in the object intent. Unfortunately, this simple observation does not allow to correct all the errors in implicative theory.

**Example 2.** Consider Case4 from Fig. 5. Case4 has set of attributes  $A = \{\text{has equal legs, has equal angles, all legs equal, at least 3 different legs}\}$ . The closure  $A''$  in the context from Fig. 4 is equal to the set of all attributes  $M$ . Therefore, closure approach would ask if the user has forgotten to add all the attributes that are still possible to add. The suggestion to add all other attributes is not supported by any example in the context as there are no objects with all attributes. More than that, such solution is not minimal in general. Therefore, such a solution is not satisfactory.

A more general description of the situation in the example above is the following. Let  $A \subseteq M$  be the intent of the inspected object such that  $\nexists g \in G : A \subseteq g'$ . In this case  $A'' = M$  and the implication  $A \rightarrow A'' \setminus A$  has an empty support. We could try to solve this problem by allowing to ask only those questions that have a supporting example in the context.

**Example 3.** Consider again Case4 from Fig. 5. As support for every question is required only the following question would be asked: has equal legs, has equal angles, at least 3 different legs  $\rightarrow$  at least 3 different angles? Support: {Quadrangle with 2 equal legs and 2 equal angles, rectangular trapezium with 2 equal legs}. However, a smaller and more intuitive correction would be to suggest the user to remove the attribute “at least 3 different legs”. If this is indeed the source of error then even after adding the suggested attribute the error would not be eliminated and would impact the implicative theory.

At this point we conclude that it is necessary to be able to suggest corrections for errors of Class 2. Such errors may be present if the object intent contains subset of attributes that none of the objects in the context has.

#### 4.1. Crucial implications

Suppose we have a new object  $g_n$  with intent  $A$  and we want to see whether  $A$  respects (is consistent with) the previous knowledge given by the context, i.e., does not have errors of Class 1 or 2. In order to find errors of Class 1 we need to know, whether, according to implications (implicative dependencies) of the context, the new object should have more attributes than just  $A$ . If this is the case, there should be an implication  $B \rightarrow c$  not respected by  $A$ :  $B \subseteq A$ , but  $c \notin A$ . Similarly, in order to find errors of Class 2, we look for implications  $B \rightarrow \bar{c}$  such that  $B \subseteq A$ , but  $c \in A$ . The following proposition shows that we do not need to look for all such implications, but for a much smaller subset of them.

**Proposition 1.** Let  $\mathbb{K} = (G, M, I)$ ,  $g_n = A$ ,  $A \subseteq M$ . Let

$$\mathcal{I}_A^{(\mathbb{K})} = \{B \rightarrow c \mid B \in \mathcal{MC}_A, c \in (B'' \setminus A) \cup (\overline{A \setminus B})\},$$

where  $\mathcal{MC}_A = \{B \in \mathcal{C}_A \mid \nexists C \in \mathcal{C}_A : B \subset C\}$  and  $\mathcal{C}_A = \{A \cap g' \mid g \in G\}$ . The set  $\mathcal{I}_A^{(\mathbb{K})}$  contains (unit) implications with nonempty support that are valid in  $\mathbb{K}$  and not respected by  $A$ . If an implication  $(E \rightarrow d)$ ,  $E \subseteq A$ ,  $d \in (M \setminus A) \cup \bar{A}$ , with nonempty support is valid in  $\mathbb{K}$ , then there is an implication  $(B \rightarrow d) \in \mathcal{I}_A^{(\mathbb{K})}$  such that  $E \subseteq B \subseteq A$ .

The last statement says that the set of implications  $\mathcal{I}_A^{(\mathbb{K})}$  is enough to deduce every attribute that can be deduced from implications of the context with nonempty support. Implications from  $\mathcal{I}_A^{(\mathbb{K})}$  are called  $A$ -crucial in  $\mathbb{K}$ . If ambiguity is excluded we omit the upper index and write simply  $\mathcal{I}_A$ .

**Proof.** Let  $(B \rightarrow c) \in \mathcal{I}_A$ , hence  $B' \subseteq c'$  by the definition of implication. By the definition of  $\mathcal{I}_A$  one has  $B = A \cap g'$  for some  $g \in G$ . Then  $B \subseteq g'$  and by the antimonotonicity of  $(\cdot)'$  one has  $g'' \subseteq B'$ . Hence,  $g'' \subseteq B' \subseteq c'$  and  $c'' \subseteq g''' = g'$ . Since  $c \in c''$ , one has  $c \in g'$ . Since  $c \in g'$  and  $B' \subseteq g$ , by the properties of  $(\cdot)'$ , one has  $(B \cup c)' = B' \cap c' \subseteq g$ . Hence, the support of  $B \rightarrow c$  contains  $g$  and is not empty. Consider the following possible cases:

1.  $c \in B'' \setminus A$ . Since  $B'' \setminus A \not\subseteq A$  implication  $B \rightarrow c$  is not respected by  $A$ ;
2.  $c \in \overline{A \setminus B}$ . Since  $A \setminus B \subseteq A$ , one has  $c \notin A$  and  $B \rightarrow c$  is not respected by  $A$ .

Now let  $E \rightarrow d$  be a valid implication not respected by  $A$  with a nonempty support. Then  $E \subseteq A$ ,  $d \notin A$  and there exists  $g \in G$  such that  $E \subseteq g'$ ,  $d \in g'$ . Therefore, there exists  $B_{\mathcal{C}_A} \in \mathcal{C}_A$  such that  $B_{\mathcal{C}_A} = A \cap g'$  and  $E \subseteq B_{\mathcal{C}_A}$ . Moreover, there exists  $B_{\mathcal{MC}_A} \in \mathcal{MC}_A$  such that  $B_{\mathcal{C}_A} \subseteq B_{\mathcal{MC}_A}$ . By construction  $B_{\mathcal{MC}_A} \subseteq A$ , therefore,  $E \subseteq B_{\mathcal{MC}_A} \subseteq A$ . Consider the following possible cases:

1.  $d \in M$ . As  $E \subseteq B_{\mathcal{MC}_A}$  by the properties of  $(\cdot)'$  one has  $E'' \subseteq B''_{\mathcal{MC}_A}$ . By the definition of the validity of an implication one has  $d \in E''$ , hence,  $d \in B''_{\mathcal{MC}_A}$ . Therefore,  $(B_{\mathcal{MC}_A} \rightarrow d) \in \mathcal{I}_A$ ;
2.  $d \in \bar{M}$ . Let  $\bar{c} = d$ . For any  $B \in \mathcal{MC}_A$  there exists  $g_* \in G$  such that  $B = A \cap g'_*$ . If  $E \subseteq g'_*$  then by the validity of the implication one has  $d \in g'_*$ , hence,  $c \notin g'_*$ . Therefore,  $c \notin B$ . As  $d \notin A$  one has  $c \in A$ . Hence,  $c \in A \setminus B$  and  $d \in \overline{A \setminus B}$ . Therefore,  $(B_{\mathcal{MC}_A} \rightarrow d) \in \mathcal{I}_A$ .  $\square$

**Proposition 2.** For a new object  $g$  with intent  $A$  one has  $\mathcal{I}_A \leq |G| \times |M|$ .

**Proof.** By the definition of  $\mathcal{MC}_A$  it contains no more than  $|G|$  elements. For any  $A, B \subseteq M$  one has  $|B''| \leq |M|$ , hence  $|B'' \setminus A| \leq |M| - |A|$ ;  $|\overline{A \setminus B}| \leq |A|$ . Hence,  $|(B'' \setminus A) \cup (\overline{A \setminus B})| \leq |(B'' \setminus A)| + |\overline{A \setminus B}| \leq |M| - |A| + |A| = |M|$ . Therefore,  $\mathcal{I}_A$  contains not more than  $|G| \times |M|$  implications.  $\square$

According to Proposition 2 to check errors of Classes 1 and 2 one has to consider polynomially many implications instead of exponentially many implications in the cardinality-minimum canonical base [20].

Proposition 1 allows one to design an algorithm for computing the set of questions (in form of implications) that can help to reveal possible errors of Classes 1 and 2.

**Proposition 3.** Let  $g$  be a new object with intent  $A$ .  $\mathcal{I}_A$  can be computed in  $O(|G|^2 \times |M|)$  time.

**Proof.** Consider the following `inspect_closure` algorithm

```

Input:  $\mathbb{K} = (G, M, I)$ ,  $A \subseteq M$ 
Output:  $\mathcal{I}_A$ 
1 if  $A'' = A$  then
2   return  $\emptyset$ 
3 Candidates = {object'  $\cap$  A | object  $\in$  G}
4 MaxCandidates = {C  $\in$  Candidates |  $\nexists B \in$  Candidates:  $C \subseteq B$ }
5 Result =  $\emptyset$ 
6 for Candidate in MaxCandidates do
7   Result.add(({Candidate  $\rightarrow$  d | d  $\in$  (Candidate''  $\setminus$   $\overline{A \cup \bar{A} \setminus$  Candidate)}))
8 return Result

```

Here  $A$  is the intent of the new object. In line 3 the algorithm computes the set of all subsets that are candidates for the premises of crucial implications. In line 4 all non-maximal subsets are discarded. In lines 6 and 7 closures of the premises are computed and the corresponding implications are added to the set of crucial implications. To estimate the worst-case complexity of the algorithm, note that executing line 1 and line 3 take at most  $O(|G| \times |M|)$  time, line 4 takes  $O(|M|)$  time for each of  $O(|G|^2)$  containment tests, and lines 6 and 7 take  $O(|G| \times |M|)$  time for computing closures of at most  $O(|G|)$  premises of crucial implications. Hence, the total worst-case time complexity is  $O(|G|^2 \times |M|)$ .  $\square$

**Example 4.** Consider the context  $\mathbb{K}_m$  in Fig. 2. We compute  $\mathcal{I}_{Case2^e}^{(\mathbb{K}_m)}$  for Case2 from the context  $\mathbb{K}_m$  in Fig. 3. In order to compute  $\mathcal{MC}_{Case2^e}$  we first compute  $\mathcal{C}_{Case2^e}$ . In order to do this we intersect  $Case2^e$  with intent of every object in the context. For example,  $Case2^e \cap (RT)' = \{\text{has right angle}\}$ . However,  $Case2^e \cap (RT \text{ with 2 equal legs})' = \{\text{has equal legs, has right angle}\} \supseteq \{\text{has right angle}\}$ , hence,  $\{\text{has right angle}\}$  will not be in  $\mathcal{MC}_{Case2^e}$ . As  $\{\text{all angles equal, all legs equal}\} \subseteq Case2^e$  and  $\{\text{all angles equal, all legs equal}\} \not\subseteq (RT)'$  we obtain the first implication from (1). As  $Case2^e$  and  $RT'$  differ in no more than  $|M|$  attributes the length of every implication is not larger than  $|M|$ . Therefore, we have at most  $|G|$  implications of size at most  $|M|$ . If we convert the implications to unit form the size of this set of implication will be limited by  $|G| \times |M|$ . Finally,

$$\mathcal{I}_{Case2^e} = \left\{ \begin{array}{l} \text{has right angle, has equal legs} \rightarrow \overline{\text{all angles equal, all legs equal}} \\ \text{has right angle, has equal legs} \rightarrow \overline{\text{at least 3 different legs,}} \\ \text{at least 3 different angles} \end{array} \right\} \quad (1)$$

The first implication is supported by Q with 2 equal legs and right angle and RT with 2 equal legs and not violated by any object from  $\mathbb{K}_m$ . The second implication is supported by the same object and is also valid in  $\mathbb{K}_m$ . Both implications are not respected by Case2.

If there are several new objects, the set of crucial implications for each new object is in general dependent on the order of adding objects.

The following statement shows which additional questions should be asked in order to compensate for this dependency.

**Proposition 4.** Let  $g_1$  and  $g_2$  be new objects with intents  $A_1$  and  $A_2$ , respectively. Let  $\mathbb{K}_1 = (G \cup g_1, M, I \cup \{(g_1, m) \mid m \in A_1\})$ . If  $\nexists g \in G : A_1 \cap A_2 \subseteq g'$  then  $I_{A_2}^{(\mathbb{K}_1)} \setminus I_{A_2}^{(\mathbb{K})} = \{A_2 \cap A_1 \rightarrow m \mid m \in (A_1 \setminus A_2) \cup (\overline{A_2 \setminus A_1})\}$ , otherwise  $I_{A_2}^{(\mathbb{K}_1)} = I_{A_2}^{(\mathbb{K})}$ .

Convex quadrangles	has equal legs	has equal angles	has right angle	all legs equal	all angles equal	at least 3 different angles	at least 3 different legs
Quadrangle (Q)						×	×
Parallelogram	×	×					
Rectangular trapezium (RT)		×	×			×	×
Q with 2 equal legs and right angle	×		×			×	×
Isosceles trapezium	×	×				×	
RT with 2 equal legs	×	×	×			×	×
Q with 2 equal angles		×				×	×
Q with 2 equal legs	×					×	×
Q with 2 equal legs and 2 equal angles	×	×				×	×

Fig. 2. Minified context of convex quadrangles  $\mathbb{K}_m$ .

Tentative errors	has equal legs	has equal angles	has right angle	all legs equal	all angles equal	at least 3 different angles	at least 3 different legs
Square	×	×	×	×	×		
Case2	×		×	×	×		

Fig. 3. Minified context of new quadrangles  $\mathbb{K}_{em}$ .

**Proof.** By the definition of  $\mathcal{I}_A^{(\mathbb{K})}$  only maximal intersections may become premises of implications. Hence, if there exists  $g \in G$  such that  $A_1 \cap A_2 \subseteq g'$  then no new implications can arise. However, if such  $g$  does not exist, the set of new implications, by the definition of  $\mathcal{I}_A^{(\mathbb{K})}$ , is  $\{(A_1 \cap A_2) \rightarrow m \mid m \in (A_1 \cap A_2)'' \setminus A_2 \cup \overline{A_2 \setminus (A_1 \cap A_2)}\}$ . As  $A_2 \setminus (A_1 \cap A_2) = A_2 \setminus A_1$  and, by assumption about maximal intersection,  $(A_1 \cap A_2)'' = A_1$ , the set of new implications is  $\{A_2 \cap A_1 \rightarrow m \mid m \in (A_1 \setminus A_2) \cup \overline{(A_2 \setminus A_1)}\}$ .  $\square$

**Example 5.** Consider contexts  $\mathbb{K}_m$  from Fig. 2 and  $\mathbb{K}_{em}$  from Fig. 3. If first square is added to  $\mathbb{K}_m$  and then Case2 is checked for errors, then the implication

has right angle, has equal legs, all legs equal, all angles equal  $\rightarrow$  has equal angles

will be output to the user. However, if Case2 is checked before square is added, the implication above will not be asked as there are no objects having the attributes “all legs equal” or “all angles equal” in the context.

Obviously, no more than  $|(A_1 \setminus A_2) \cup (A_2 \setminus A_1)|$  additional questions may arise. It is also easy to see that additional implications only arise in case where both new objects have maximal intents in  $\mathbb{K}$ . However, as we do not require any information about maximal intents in data domain, we have to be careful when adding an object with maximal intent. The following corollary states clearly which implications should be considered in order to guarantee the absence of errors that may affect the implication base.

**Proposition 5.** Let  $g_n$  be a new object with intent  $A$ . The sets of implications  $I_1 = \{A \rightarrow m \mid m \in M \setminus A, \nexists g \in G : A \subseteq g'\}$  and  $I_2 = \{(A - a) \rightarrow \bar{a} \mid a \in A, \nexists g \in G : (A - a) \subseteq g'\}$  are valid in  $\mathbb{K}$ , have empty support, and are not respected by  $A$ .

Note that if there exists  $g \in G$  such that  $A \subseteq g'$  then both sets  $I_1$  and  $I_2$  are empty.

Convex quadrangles	has equal legs	has equal angles	has right angle	all legs equal	all angles equal	at least 3 different angles	at least 3 different legs
Square	×	×	×	×	×		
Rectangle	×	×	×		×		
Quadrangle (Q)						×	×
Rhombus	×	×		×			
Parallelogram	×	×					
Rectangular trapezium (RT)		×	×			×	×
Q with 2 equal legs and right angle	×		×			×	×
Isosceles trapezium	×	×				×	×
RT with 2 equal legs	×	×	×			×	×
Q with 2 equal angles		×				×	×
Q with 2 equal legs	×					×	×
Q with 2 equal legs and 2 equal angles	×	×				×	×

Fig. 4. Context of convex quadrangles  $\mathbb{K}$ .

Tentative errors	has equal legs	has equal angles	has right angle	all legs equal	all angles equal	at least 3 different angles	at least 3 different legs
Case1	×			×		×	
Case2	×		×	×	×		
Case3		×	×	×	×	×	×
Case4	×	×		×			×

Fig. 5. Context of tentative errors  $\mathbb{K}_e$ .

**Proof.** By definition  $B \rightarrow c$  is respected by  $N \subseteq M$  if  $B \not\subseteq N$ . By the definition of  $I_1$  and  $I_2$  all the premises are not contained in any object intents from the context. Therefore, implications are valid, however, they are not supported by any of the object intents.

As  $A \subseteq A$ ,  $(M \setminus A) \not\subseteq A$ , implications from  $I_1$  are not respected by  $A$ . As  $\forall a \in A : (A - a) \subseteq A$ ,  $a \notin (M \setminus A)$ , implications from  $I_2$  are not respected by  $A$ .  $\square$

According to Proposition 5 the number of additional questions for new objects that have maximal intents cannot exceed  $|M|$ . As none of the questions have objects from context in support we suggest that maximal objects should be checked “by hand”.

For the sake of compactness in what follows we present implications in non-unit form. The name `inspect_base` is used to denote the function implementing base approach.

#### 4.2. Example

Consider the following example with convex quadrangles. Formal context given by the cross-table in Fig. 4 contains convex quadrangles and their properties. The context does not cover the domain completely, i.e. not all possible convex quadrangle types are considered. Attributes “has equal legs” and “has equal angles” require at least two angles/legs of a quadrangle to be equal. Some dependencies on attributes are trivial, e.g., if all angles in a quadrangle are equal, then this quadrangle has equal angles.



Four objects are in the context of tentative errors in Fig. 5. These objects are added to the context in Fig. 4 one at a time.

Inspecting Case1:

```
inspect_base
  all legs equal → has equal angles, has equal legs
inspect_closure
  has equal legs, at least 3 different angles → at least 3 different legs, all legs equal
  has equal legs, all legs equal → has equal angles, at least 3 different angles
```

Both algorithms reveal possible errors in a similar manner, although there are obvious differences. In the output of `inspect_base` the premises are smaller than in the output of `inspect_closure`. The latter also reveals dependencies of Class 2. It is easy to see that all output implications hold in data domain. For example, if all legs are equal in a quadrangle, it should have equal angles and should not have 3 different angles. Hence, this object should be recognized as an error, it should be corrected to rhombus or to quadrangle with two equal legs.

Inspecting Case2:

```
inspect_base
  all angles equal → has equal angles, has equal legs, has right angle
  all legs equal → has equal angles, has equal legs
inspect_closure
  has right angle, has equal legs, all legs equal, all angles equal → has equal angles
```

In this example we are able to ask even less questions to an expert using `inspect_closure` as with `inspect_base`. This is the result of finding implications generated by maximal subsets of object's intent. The intent of Case2 occurs in the context (in the intent of Square), that is why we do not get any negative attributes in the output of `inspect_closure`. Again, all implications are valid in data domain, therefore, Case2 is an error.

Inspecting Case3:

```
inspect_base
  all angles equal → has equal angles, has equal legs, has right angle
  all legs equal → has equal angles, has equal legs
inspect_closure
  has equal angles, has right angle, at least 3 different legs, at least 3 different angles → all angles equal, all legs equal
  has equal angles, has right angle, all legs equal, all angles equal → has equal legs, at least 3 different angles,
  at least 3 different legs
```

In Case3 we get both implications from the output of `inspect_base` combined in one implication with a bigger premise in the output of `inspect_closure`. In addition we obtain several implications with negative attributes. It is easy to see that all implications hold in the data domain, therefore, Case3 is an error and should be corrected either to rectangular trapezium or to square.

Inspecting Case4:

```
inspect_base
  has equal angles, has equal legs, at least 3 different legs, all legs equal → has right angle, at least 3 different angles,
  all angles equal
inspect_closure
  has equal angles, has equal legs, all legs equal → at least 3 different legs
  has equal angles, has equal legs, at least 3 different legs → at least 3 different angles, all legs equal
```

Case4 is a very special case where the corresponding implication from canonical base has empty support. In the output of `inspect_base` we obtain all questions possible for this intent. As discussed above these questions are not based on any information input so far. The reason for that is that Case4 has maximal intent in the context. So these questions could also be found using Proposition 5. However, even if we add attributes "at least 3 different angles" and "all angles equal" and reject the last implication we would not be able to recognize this object as an error. On the contrary `inspect_closure` allows us to recognize errors of Class 2 and state that Case4 should be corrected to have the intent of rhombus or quadrangle with two equal legs and two equal angles.

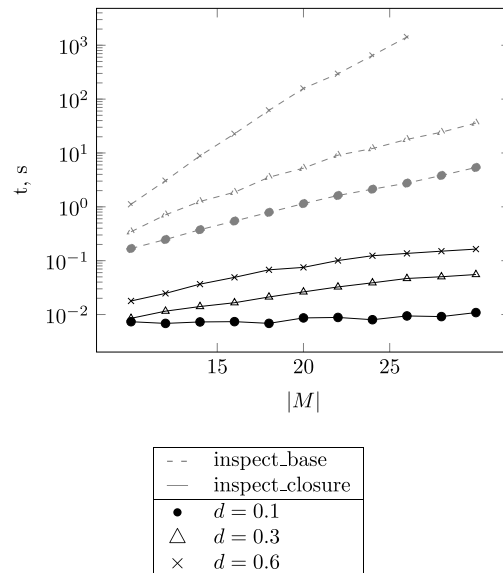


Fig. 6. Comparison of runtime on synthetic contexts in semilog scale.

## 5. Experiment

Below the results of experiments on synthetic data are presented. The experiments were conducted as follows: all objects are first taken one by one out from the context and then added as new objects; all the possible errors of Classes 1 and 2 are found and output. In this experiment we wanted to compare the runtime of the algorithms.

An FCA package for Python was used for implementation [1]. For computing the canonical base an optimized algorithm based on Next Closure was used [25]. All tests described below were run on computer with Intel Core i7 1.6 GHz processor and 4 Gb of RAM running Linux Ubuntu 11.10 x64.

In Fig. 6 the results of running both algorithms on synthetic contexts are presented. For each context the number of objects is equal to 50. Parameter  $d$  represents the density of the context, i.e. the probability of having a cross in the cross-table representing the relation. This result is presented in the semi-logarithmic scale. It is easy to note that with the growth of the number of attributes and the density, the difference between runtime of two algorithms grows as well.

Another experiment was conducted to test the quality of finding errors by the introduced method. The information about dependencies between negative attributes is not reflected in the implication base. Therefore, more implications are usually violated by objects having more attributes in their intent. Small intents usually violate only few implications. However, in this experiment we aim at finding not only the errors affecting the implication base; therefore, it is necessary to level out the shift between larger and smaller intents. For this purpose in the following experiments a slight modification of the introduced method is used. The complementary context to a context  $\mathbb{K} = (G, M, I)$  is defined as  $\mathbb{K}^c := (G, M, (G \times M) \setminus I)$ . The method applied to the complementary context will output implications with only negative attributes in the premise. Running the introduced method on both original context and complementary context yields better results. Note that implications with both positive and negative attributes will not be generated.

The experiments were conducted in the following settings. An object was picked up from a context, from one to three errors were randomly introduced into the intent of the object. The method was used to find possible errors in the object. If all the erroneous attributes were in the conclusion of the unit implications with the same premise then the errors were marked as found. In this case all the erroneous attributes are contained in one question to the user. Afterwards the object already without errors is returned to the context and the next object is picked up. We considered three contexts from the UCI repository [11]: SPECT, house-votes-84 and kr-vs-kp. Therefore, nine experiments were conducted. In every experiment 1000 objects (with possible repetitions) were picked up one after another. In Table 1 the results of the experiment are presented.

The more objects there are in the context, the less is the number of all valid implications of the context (which does not mean the smaller size of the implication base). Therefore, less invalid implications could be output to the user. In SPECT intents are very diverse (there are only 5.78 irreducible objects per attribute on average) that is why not more than three implications on average are output and bad ratio of found errors is obtained. In house-votes-84 intents are more similar, that is why we have more questions per object. The ratio of found errors for one error is relatively high, however, it quickly drops with the increase of errors, as the number of irreducible objects is small. In kr-vs-kp there are much more objects per attribute and the results for two and three errors at a time are much better. However, if the user is able to correctly answer all unit implications even better results can be achieved. In this case the user may correct first errors and repeat

**Table 1**

Error finding experiment carried out on three contexts from UCI.

Context name	G	G  after reducing	M
SPECT	267	133	23
house-votes-84	232 <sup>a</sup>	104	17
kr-vs-kp	3198	453	36 <sup>b</sup>

Number of errors per object	Errors found	Found/all ratio	Total number of implications	Implications per object
		(a) SPECT		
1	548	0.548	2298	2.41
2	242	0.242	2753	2.80
3	131	0.131	2703	2.71
		(b) house-votes-84		
1	712	0.712	9780	11.4
2	217	0.217	14018	14.7
3	71	0.071	18276	18.4
		(c) kr-vs-kp		
1	786	0.786	7520	8.47
2	393	0.393	12863	13.2
3	247	0.247	18322	18.3

<sup>a</sup> All objects containing missing values were removed.<sup>b</sup> Attribute 15 was removed due to many-valuedness.

the procedure having already only one error. In these experiments the error-finding process was considered successful only if there is one implication suggesting all the needed corrections at once.

It is worth noting that the chance of random guess in predicting all errors in an object description is only  $1/(|M|^n)$  if there are  $n$  errors as compared to 50% for the classification task.

The results testify to the obvious advantage of applying `inspect_closure` approach to `inspect_base` approach in terms of runtime.

## 6. Conclusion

A method for finding errors in implicative theories was introduced. The method uses some techniques based on Formal Concept Analysis. As opposed to finding the canonical (cardinality minimal) base of implications, which can be very time consuming due to intrinsic intractability, the proposed algorithm terminates in polynomial time. Moreover, after checking maximal object descriptions (object intents) “by hand” it is possible to find all errors of two considered types or prove their absence. Computer experiments show that in practice the proposed method works much faster than that based on the generation of the implication base.

## Acknowledgements

The first author was supported by the Basic Research Program of the National Research University Higher School of Economics, project “Mathematical models, algorithms, and software for knowledge discovery in structured and text data”. The second author was supported by German Academic Exchange Service (DAAD). We thank Bernhard Ganter and Sergei Obiedkov for discussion and useful remarks.

## References

- [1] N. Romashkin, A. Revenko, Python package for formal concept analysis, <https://github.com/artreven/fca>.
- [2] F. Baader, D. Calvanese, D.L. McGuinness, D. Nardi, P.F. Patel-Schneider, *The Description Logic Handbook: Theory, Implementation and Applications*, 2010.
- [3] F. Baader, R. Penaloza, B. Suntisiraparorn, Pinpointing in the description logic  $\mathcal{EL}^+$ , in: *KI 2007: Advances in Artificial Intelligence*, Springer, 2007, pp. 52–67.
- [4] M.A. Babin, S.O. Kuznetsov, Computing premises of a minimal cover of functional dependencies is intractable, *Discrete Appl. Math.* 161 (6) (2013) 742–749.
- [5] V.G. Blinova, D.A. Dobrynin, V.K. Finn, S.O. Kuznetsov, E.S. Pankratova, Toxicology analysis by means of JSM-method, *Bioinformatics* 19 (2003) 1201–1207.
- [6] A. Buzmakov, E. Egho, N. Jay, S.O. Kuznetsov, A. Napoli, C. Raissi, On projections of sequential pattern structures (with an application on care trajectories), in: Jan Outrata, Manuel Ojeda-Aciego (Eds.), *10th International Conference on Concept Lattices and Their Applications, CLA 2013*, 2013, pp. 199–208.
- [7] F. Dau, Implications of properties concerning complementation in finite lattices, in: D. Dorninger, et al. (Eds.), *Proceedings of the 58th Workshop on General Algebra “58. Arbeitstagung Allgemeine Algebra”, Vienna, Austria, June 3–6, 1999*, in: *Contrib. Gen. Algebra*, vol. 12, Verlag Johannes Heyn, Klagenfurt, 2000, pp. 145–154.

- [8] F. Distel, B. Sertkaya, On the complexity of enumerating pseudo-intents, *Discrete Appl. Math.* 159 (6) (2011) 450–466.
- [9] W. Fan, J. Li, S. Ma, N. Tang, W. Yu, Towards certain fixes with editing rules and master data, *Proc. VLDB Endow.* 3 (1) (2010) 173–184.
- [10] F. Baader, B. Sertkaya, Usability issues in description logic knowledge base completion, in: S. Rudolph, S. Ferré (Eds.), 7th International Conference on Formal Concept Analysis, ICFA 2009, in: *LNAI*, vol. 5548, 2009, pp. 1–21.
- [11] A. Frank, A. Asuncion, *UCI machine learning repository*, <http://archive.ics.uci.edu/ml>, 2010.
- [12] W.J. Frawley, G. Piatetsky-Shapiro, C.J. Matheus, Knowledge discovery in databases: an overview, *AI Mag.* 13 (3) (1992) 57.
- [13] B. Ganter, Two basic algorithms in concept analysis, in: L. Kwuida, B. Sertkaya (Eds.), *Formal Concept Analysis*, in: *Lect. Notes Comput. Sci.*, vol. 5986, Springer, Berlin, Heidelberg, 2010, pp. 312–340.
- [14] B. Ganter, R. Wille, *Formal Concept Analysis: Mathematical Foundations*, Springer, 1999.
- [15] J.-L. Guigues, V. Duquenne, Familles minimales d'implications informatives résultant d'un tableau de données binaires, *Math. Sci. Hum.* 24 (95) (1986) 5–18.
- [16] Y.K. Jain, V. Suryawanshi, A new approach for handling null values in web log using KNN and tabu search KNN, *Int. J. Data Min. Knowl. Manage. Process* 1 (5) (September 2011) 9–19.
- [17] H.A. Kautz, M.J. Kearns, Bart Selman, Reasoning with characteristic models, in: *The Eleventh National Conference on Artificial Intelligence, AAAI-93*, 1993, pp. 1–14.
- [18] M. Kaytoue, S.O. Kuznetsov, A. Napoli, S. Duplessis, Mining gene expression data with pattern structures in formal concept analysis, *Inf. Sci.* 181 (2011) 1989–2001.
- [19] M. Kirchberg, E. Leonardi, Y.S. Tan, S. Link, R.K.L. Ko, B.-S. Lee, Formal concept discovery in semantic web data, in: *10th International Conference on Formal Concept Analysis, ICFA 2012*, in: *LNAI*, vol. 7278, 2012, pp. 164–179.
- [20] S.O. Kuznetsov, On the intractability of computing the Duquenne–Guigues base, *J. Univers. Comput. Sci.* 10 (8) (Aug. 2004) 927–933.
- [21] S.O. Kuznetsov, S. Obiedkov, Some decision and counting problems of the Duquenne–Guigues basis of implications, *Discrete Appl. Math.* 156 (11) (2008) 1994–2003.
- [22] L. Kwuida, C. Pech, H. Reppe, Generalizations of boolean algebras. An attribute exploration, *Math. Slovaca* 56 (2) (2006) 145–165.
- [23] T. Meyer, K. Lee, R. Booth, J.Z. Pan, Finding maximally satisfiable terminologies for the description logic ALC, in: *Proceedings of the National Conference on Artificial Intelligence*, vol. 21, AAAI Press/MIT Press, Menlo Park, CA/Cambridge, MA/London, 2006, p. 269.
- [24] B. Mirkin, *Mathematical Classification and Clustering*, Kluwer Academic Publisher, 1996.
- [25] S. Obiedkov, V. Duquenne, Attribute-incremental construction of the canonical implication basis, *Ann. Math. Artif. Intell.* 49 (1–4) (April 2007) 77–99.
- [26] J. Rasmussen, The role of hierarchical knowledge representation in decisionmaking and system management, *IEEE Trans. Syst. Man Cybern.* 2 (1985) 234–243.
- [27] A. Revenko, Automated construction of implicative theory of algebraic identities of size up to 5, in: Cynthia Vera Glodeanu, Mehdi Kaytoue, Christian Sacarea (Eds.), *Formal Concept Analysis*, in: *Lect. Notes Comput. Sci.*, vol. 8478, Springer International Publishing, 2014, pp. 188–202.
- [28] A. Revenko, S.O. Kuznetsov, Finding errors in new object intents, in: *CLA 2012*, CEUR, 2012, pp. 151–162.
- [29] J.M. Rodríguez-Jimenez, P. Cordero, M. Enciso, A. Mora, A generalized framework to consider positive and negative attributes in formal concept analysis, in: *CLA 2014*, 2014.
- [30] U. Rysseel, F. Distel, D. Borchmann, Fast computation of proper premises, in: Amedeo Napoli, Vilem Vychodil (Eds.), *International Conference on Concept Lattices and Their Applications*, INRIA Nancy – Grand Est and LORIA, 2011, pp. 101–113.
- [31] E.-L. Silva-Ramírez, R. Pino-Mejías, M. López-Coello, M.-D. Cubiles-de-la Vega, Missing value imputation on missing completely at random data using multilayer perceptrons, *Neural Netw.* 24 (1) (January 2011) 121–129.
- [32] Y.L. Simmhan, D. Gannon, B. Plale, A survey of data provenance techniques, Technical report iub-cs-tr618, Computer Science Department, Indiana University, Bloomington, 2005.
- [33] Q. Song, M. Shepperd, A new imputation method for small software project data sets, *J. Syst. Softw.* (2007) 1–24.
- [34] P. Valtchev, R. Missaoui, R. Godin, Formal concept analysis for knowledge discovery and data mining: the new challenges, in: *Concept Lattices*, Springer, 2004, pp. 352–371.