

Generic association rule bases: Are they so succinct?

T. Hamrouni¹, S. Ben Yahia¹ and E. Mephu Nguifo²

¹ Computer Science department, Faculty of Sciences of Tunis, Tunis, Tunisia.

{tarek.hamrouni, sadok.benyahia}@fst.rnu.tn

² CRIL-CNRS, IUT de Lens, Lens, France.

mephu@cril.univ-artois.fr

Abstract. In knowledge mining, current trend is witnessing the emergence of a growing number of works towards defining “concise and lossless” representations. One main motivation behind is: tagging a unified framework for drastically reducing large sized sets of association rules. In this context, generic bases of association rules – whose backbone is the conjunction of the concepts of minimal generator and closed itemset (CI) – constituted so far irreducible compact nuclei of association rules. However, the inherent absence of a unique minimal generator (MG) associated to a given CI offers an “ideal” gap towards a tougher redundancy removal even from generic bases of association rules. In this paper, we adopt the succinct system of minimal generators (SSMG), newly redefined in [1], to be an *exact* representation of the MG set. Then, we incorporate the SSMG into the framework of generic bases to only maintain the *succinct* generic association rules. After that, we give a thorough formal study of the related inference mechanisms allowing to derive *all redundant* association rules starting from succinct ones. Finally, an experimental study shows that our approach makes it possible to eliminate without information loss an important number of *redundant* generic association rules and thus, to only present *succinct* and *informative* ones to users.

1 Introduction

As an important topic in data mining, association rule mining research [2] has progressed in various directions. Unfortunately, one problem with the current trend is that it mainly favoured the efficient extraction of interesting itemsets regardless the effectiveness of the mined knowledge. Indeed, by laying stress on the “algorithmic” improvement of the frequent (closed) itemset extraction step, the current trend neglects user’s needs: “concise with add-value knowledge”. Hence, the number of association rules, which can be extracted even from small datasets, is always a real hampering towards their effective exploitation by the users. Indeed, at the end of the extraction process, the user is faced to an overwhelming quantity of association rules among which a large number is *redundant*, what badly affects the quality of their interpretability. Nevertheless, some approaches have been devoted to the reduction of the number of association rules such as generic bases [3–7], concise representations [8–10], quality measures [11], user-defined templates or constraints [12, 13]. Among them, generic bases constitute an interesting starting point to reduce without loss of information the size of the association rule set. Indeed, using the mathematical settings of the Formal Concept Analysis

(FCA) [14], generic bases were flagged as irreducible nuclei of association rules from which *redundant* ones can be derived without any loss of information [3]. In this context, different works have shown that generic bases, containing association rules whose implications are between minimal generators (MGs) [3] and closed itemsets (CIs) [8], convey the maximum of information since they are of minimal premises and of maximal conclusions [3, 15]. For these reasons, such association rules are considered as the most informative ones [3].

Nevertheless, a recent study proposed by Dong *et al.* shows that the MG set still present a kind of redundancy [16]. Indeed, they consider the set of the MGs associated to a given CI by distinguishing two distinct classes: *succinct* MGs and *redundant* ones. Thus, Dong *et al.* introduce the succinct system of minimal generators (SSMG) as a concise representation of the MG set. They state that *redundant* MGs can be withdrawn from the MG set since they can straightforwardly be inferred, without loss of information, using the knowledge gleaned from the *succinct* ones [16]. However, in [1], we showed that the *succinct* MGs, as defined by Dong *et al.*, prove not to be an *exact* representation (no loss of information *w.r.t.* *redundant* MGs) in contrary to authors' claims. We also presented new definitions allowing to overcome the limitations of their work and, hence, to make of the SSMG really an *exact* representation.

In this paper, we propose to incorporate the SSMG, as redefined in [1], into the framework of generic bases to reduce as far as possible the redundancy within generic association rules. Thus, after a study of the best known generic bases of association rules, we apply the SSMG to the couple of generic bases proposed by Bastide *et al.* [3]. This couple presents at least two complementary advantages. On the one hand, association rules composing it are of minimum premises and of maximal conclusions, and, hence, convey the maximum of information [3, 15]. On the other hand, this couple gathers the *ideal* properties of an association rule representation since it is lossless, sound and informative [5]. We then study the obtained generic bases - once the SSMG is applied - to check whether they are extracted without loss of information. Finally, an experimental evaluation illustrates the potential of our approach towards offering to users a redundancy-free set of generic association rules. Please note that it is out of the scope of this paper to discuss how the *succinct* generic association rules are efficiently discovered.

The organization of the paper is as follows: Section 2 recalls some preliminary notions that will be used in the remainder of the paper. We devote Section 3 to the presentation of the main definition of the SSMG proposed in [1]. Section 4 is dedicated to the presentation of the *succinct* generic bases of association rules. In order to derive *all redundant* association rules that can be extracted from a context, an axiomatic system and a study of its main properties are also provided. In Section 5, several experiments illustrate the utility of our approach followed by a summary of our contributions and avenues for future work in Section 6.

2 Preliminary definitions

In this section, we present some notions that will be of use in the following.

Definition 1. (EXTRACTION CONTEXT) *An extraction context is a triplet $\mathcal{K} = (\mathcal{O}, \mathcal{I}, \mathcal{R})$, where \mathcal{O} represents a finite set of objects, \mathcal{I} is a finite set of items and \mathcal{R} is a binary (incidence) relation (i.e., $\mathcal{R} \subseteq \mathcal{O} \times \mathcal{I}$). Each couple $(o, i) \in \mathcal{R}$ expresses that the object $o \in \mathcal{O}$ contains the item $i \in \mathcal{I}$.*

The closure operator $''$ denotes the closure operator $\phi \circ \psi$ s.t. (ϕ, ψ) represents a couple of operators defined by $\psi : \mathcal{P}(\mathcal{I}) \rightarrow \mathcal{P}(\mathcal{O})$ s.t. $\psi(I) = \{o \in \mathcal{O} \mid \forall i \in I, (o, i) \in \mathcal{R}\}$ and $\phi : \mathcal{P}(\mathcal{O}) \rightarrow \mathcal{P}(\mathcal{I})$ s.t. $\phi(O) = \{i \in \mathcal{I} \mid \forall o \in O, (o, i) \in \mathcal{R}\}$ [17]. It induces an equivalence relation on the power set of items portioning it into distinct subsets called *equivalence classes* [18]. The largest element (w.r.t. the number of items) in each equivalence class is called a *closed itemset* (CI) [8] and the smallest ones are called *minimal generators* (MGs) [3]. The notions of *closed itemset* and of *minimal generator* are defined as follows:

Definition 2. (CLOSED ITEMSET) *An itemset $I \subseteq \mathcal{I}$ is said to be closed if and only if $I'' = I$ [8]. The support of I , denoted by $\text{Supp}(I)$, is equal to the number of objects in \mathcal{K} that contain I . I is said to be frequent if $\text{Supp}(I)$ is greater than or equal to a minimum support threshold, denoted minsupp .*

Definition 3. (ICEBERG CONCEPT LATTICE) *Let $\mathcal{FCI}_{\mathcal{K}}$ be the set of the frequent CIs extracted from a context \mathcal{K} . When the set $\mathcal{FCI}_{\mathcal{K}}$ is partially ordered with set inclusion, the resulting structure $(\hat{\mathcal{L}}, \subseteq)$ only preserves the Join operator [17]. This structure is called a join semi-lattice or an upper semi-lattice and is hereafter referred to as “Iceberg concept lattice” [4].*

Definition 4. (UPPER COVER) *The upper cover of a frequent CI f (denoted $\text{Cov}^u(f)$) consists of the frequent CIs that immediately cover f in the Iceberg concept lattice. The set $\text{Cov}^u(f)$ is given as follows: $\text{Cov}^u(f) = \{f_1 \in \mathcal{FCI}_{\mathcal{K}} \mid f \subset f_1 \wedge \nexists f_2 \in \mathcal{FCI}_{\mathcal{K}} \text{ s.t. } f \subset f_2 \subset f_1\}$.*

Definition 5. (MINIMAL GENERATOR) *An itemset $g \subseteq \mathcal{I}$ is said to be a minimal generator (MG) of a CI f , if and only if $g'' = f$ and $\nexists g_1 \subset g$ s.t. $g_1'' = f$ [3]. Thus, the set MG_f of the MGs associated to a CI f is: $\text{MG}_f = \{g \subseteq \mathcal{I} \mid g'' = f \wedge \nexists g_1 \subset g \text{ s.t. } g_1'' = f\}$.*

3 Succinct System of Minimal Generators

In this section, we briefly describe the main structural properties of the succinct system of minimal generators (SSMG) newly redefined in [1] to make of it an *exact* representation of the minimal generator (MG) set.

The set MG_f of the MGs associated to a given closed itemset (CI) f can be divided into different **equivalence subclasses**⁽¹⁾ thanks to a substitution process. The latter uses a substitution operator denoted *Subst*. This substitution operator is a partial one

¹ The term *equivalence subclasses* is used here instead of *equivalence classes* to avoid the confusion with the *equivalence classes* induced by the closure operator $''$.

allowing to substitute a subset of an itemset X , say Y , by another itemset Z belonging to the same equivalence class of Y (i.e., $Y'' = Z''$). This operator is then defined as follows:

Definition 6. [1] (SUBSTITUTION OPERATOR) *Let X, Y and Z be three itemsets such that $Y \subset X$ and $Y'' = Z''$. The substitution operator Subst , w.r.t. X, Y and Z , is defined as follows : $\text{Subst}(X, Y, Z) = (X \setminus Y) \cup Z$.*

It is shown in [1] that X and $\text{Subst}(X, Y, Z)$ have the same closure.

For each equivalence class \mathcal{C} (or equivalently, for each CI f), the substitution operator induces an equivalence relation on the set MG_f of the MGs of f portioning it into distinct equivalence subclasses. The definition of an equivalence subclass requires that we define the notion of *redundant* MG under the substitution process point of view as follows:

Definition 7. [1] (MINIMAL GENERATORS' REDUNDANCY) *Let g and g_1 be two MGs belonging to the same equivalence class induced by the closure operator $''$.*

- g is said to be a **direct redundant** (resp. derivable) with respect to (resp. from) g_1 , denoted $g_1 \vdash g$, if $\text{Subst}(g_1, g_2, g_3) = g$ with $g_2 \subset g_1$ and $g_3 \in \text{MG}_{\mathcal{K}}$ s.t. $g_3'' = g_2''$. The operator \vdash is reflexive, symmetric but not necessarily transitive.

- g is said to be a **transitive redundant** with respect to g_1 , denoted $g_1 \vDash g$, if it exists a sequence of n MGs ($n \geq 2$), $gen_1, gen_2, \dots, gen_n$, such that $gen_i \vdash gen_{i+1}$ ($i \in [1..(n-1)]$) with $gen_1 = g_1$ and $gen_n = g$. The operator \vDash is reflexive, symmetric and transitive.

For $n = 2$, the operator \vDash is reduced to the operator \vdash .

The definition of a *succinct* minimal generator that we give hereafter requires that we adopt a total order relation among itemsets defined as follows.

Definition 8. (TOTAL ORDER RELATION) *Let \preceq be a total order relation among item literals, i.e., $\forall i_1, i_2 \in \mathcal{I}$, we have $i_1 \preceq i_2$ or $i_2 \preceq i_1$. This relation is extended to also cope with itemsets of different sizes by first considering their cardinality. This is done as follows: Let X and Y be two itemsets and i an item s.t. $i \notin X$ and $i \notin Y$. Let $\text{Card}(X)$ and $\text{Card}(Y)$ be the respective cardinalities of X and Y . We then have:*

- $\text{Card}(X) < \text{Card}(Y) \implies X \prec Y$.
- $X \preceq Y \iff X \cup \{i\} \preceq Y \cup \{i\}$.

Example 1. If we consider the lexicographic order as the total order relation \preceq , then ⁽²⁾:

- $|d| < |be| \implies d \prec be$.
- $abd \preceq abe \iff abd \cup \{c\} \preceq abe \cup \{c\}$ (i.e., $abcd \preceq abce$).

Definition 9. [1] (EQUIVALENCE SUBCLASSES) *The operator \vDash induces an equivalence relation on the set MG_f , of the MGs associated to a CI f , portioning it into distinct subsets called equivalence subclasses. If $g \in \text{MG}_f$, then the equivalence subclass of g , denoted by $[g]$, is the subset of MG_f consisting of all elements that are*

² We use a separator-free form for the sets, e.g., be stands for $\{b, e\}$.

transitive redundant w.r.t. g . In other words, we have: $[g] = \{g_1 \in \text{MG}_f \mid g \models g_1\}$. The smallest MG in each equivalence subclass, w.r.t. the total order relation \preceq , will be considered as its **succinct** MG. While, the other MGs will be qualified as **redundant** MGs.

Example 2. Let us consider the extraction context \mathcal{K} depicted by Figure 1 (Left). The total order relation \preceq is set to the lexicographic order. Figure 1 (Right) shows, for each CI, the following information: its MGs, its *succinct* MGs and its support. The MG “ adg ” is a *succinct* one, since it is the smallest MG, w.r.t. \preceq , among those of “ $abcdeg$ ”. Indeed, when extracting the first equivalence subclass associated to “ $abcdeg$ ”, the whole MG set associated to “ $abcdeg$ ” is considered. We then have: $adg \preceq aeg$, $adg \preceq bdg$ and $adg \preceq beg$. The MG “ aeg ” is a redundant one since $\text{Subst}(adg, ad, ae) = aeg \in \text{MG}_{abcdeg}$ ($adg \vdash aeg$ and, hence, $adg \models aeg$). It is the same for the MGs “ bdg ” and “ beg ” since $adg \models bdg$ and $adg \models beg$.

	a	b	c	d	e	f	g
1			\times	\times	\times	\times	\times
2	\times	\times	\times	\times	\times		
3	\times	\times	\times			\times	\times
4	\times	\times	\times	\times			\times

	CI	MGs	Succinct MGs	Support
1	c	\emptyset	\emptyset	4
2	abc	a, b	a, b	3
3	cde	d, e	d, e	3
4	cg	g	g	3
5	cfg	f	f	2
6	$abcde$	ad, ae, bd, be	ad	2
7	$abcg$	ag, bg	ag	2
8	$abcfg$	af, bf	af	1
9	$cdeg$	dg, eg	dg	2
10	$cdefg$	df, ef	df	1
11	$abcdeg$	adg, aeg, bdg, beg	adg	1

Fig. 1. (Left) An extraction context \mathcal{K} . (Right) The CIs extracted from \mathcal{K} and for each one, the corresponding MGs, *succinct* MGs and support.

The succinct system of minimal generators (SSMG) is then defined as follows [1]:

Definition 10. [1] (SUCCINCT SYSTEM OF MINIMAL GENERATORS) A *succinct system of minimal generators (SSMG)* is a system where only *succinct* MGs are retained among all MGs associated to each CI.

Proposition 1. [1] The SSMG is an exact representation of the MG set.

In the remainder, the set of *succinct* (resp. *redundant*) frequent MGs that can be extracted from a context \mathcal{K} will be denoted $\mathcal{FMG}_{\text{suc}\mathcal{K}}$ (resp. $\mathcal{FMG}_{\text{red}\mathcal{K}}$).

4 Succinct and informative association rules

We now put the focus on integrating the concept of succinct system of minimal generators (SSMG) within the generic association rule framework. Our purpose is to obtain,

without information loss, a more compact set of all association rules, from which the remaining *redundant* ones can be generated if desired.

4.1 Association rules: some basic notations

The formalization of the association rule extraction problem was introduced by Agrawal *et al.* [2]. The derivation of association rules is achieved starting from a set of *frequent* itemsets [19] extracted from a context \mathcal{K} (denoted $\mathcal{FT}_{\mathcal{K}}$), for a minimal support threshold *minsupp*. An association rule R is a relation between itemsets and is of the form $R: X \Rightarrow (Y \setminus X)$, such that X and Y are *frequent* itemsets, and $X \subset Y$. The itemsets X and $(Y \setminus X)$ are, respectively, called the *premise* and the *conclusion* of the association rule R (also called *antecedent* and *consequent* of R [3], and *condition* and *consequence* of R [15]). An association rule is said to be *valid* (or *strong*) if its confidence measure, $\text{Conf}(R) = \frac{\text{Supp}(Y)}{\text{Supp}(X)}$, is greater than or equal to a minimal threshold of confidence denoted *minconf*. If $\text{Conf}(R) = 1$, then R is called *exact association rule*, otherwise it is called *approximate association rule*.

4.2 Extraction of succinct and informative association rules

The problem of the relevance and the usefulness of association rules is of paramount importance. Indeed, an overwhelming quantity of association rules can be extracted even from small real-life datasets, among which a large number is *redundant* (*i.e.*, conveying the same information) [4, 6]. This fact boosted the interest in novel approaches aiming to reduce this large association rule list, while preserving the most interesting rules. These approaches are mainly based on the battery of results provided by the Formal Concept Analysis (FCA) mathematical settings [14]. Thus, they focused on extracting irreducible nuclei of all association rules, commonly referred to as “*generic bases*”, from which the remaining *redundant* association rules can be derived. Definition 11 describes the properties that characterize a generic basis once it is extracted without loss of information.

Definition 11. A generic basis \mathcal{B} , associated with an appropriate inference mechanism, is said to fulfill the ideal properties of an association rule representation if it is [5]:

1. **lossless:** \mathcal{B} must enable the derivation of all valid association rules, and
2. **sound:** \mathcal{B} must forbid the derivation of association rules that are not valid, and
3. **informative:** \mathcal{B} must allow to exactly retrieve the support and confidence values of each derived association rule.

The generic basis \mathcal{B} is said to verify the property of derivability if it is lossless and sound.

The majority of the generic bases that were proposed in the literature convey association rules presenting implications between minimal generators (MGs) and closed itemsets (CIs) [3, 5, 7]. Indeed, it was proven that such association rules, with minimal premises and maximal conclusions, convey the maximum of information [3, 15] and are hence qualified as the most informative association rules [3]. Furthermore, *succinct*

MGs are very well suited for such association rules since they offer the smallest possible premises. Indeed, they are the smallest ones in their respective equivalence subclasses. They are also the most interesting ones since correlations in each *succinct* MG can not be predicted given correlations of its subsets and those of the other (redundant) MGs.

Hence, in order to extract much more compact sets of association rules, we propose to integrate the concept of the succinct system of minimal generators (SSMG) within the framework of generic bases. Although, our approach can be applied to different generic bases, we concentrate our presentation on the couple $(\mathcal{GB}, \mathcal{RI})$ of generic association rule bases proposed by Bastide *et al.* [3]. Indeed, in addition to the quality of the conveyed knowledge, the selected couple has the advantage to fulfill the *ideal* association rule representation's properties (summarized by Definition 11) in comparison to other generic bases (like the couple $(\mathcal{DGB}, \mathcal{LB})$ [4], \mathcal{RR} [5], \mathcal{NRR} [6], etc. ⁽³⁾) [5]. Moreover, as this will be shown in the continuation, these properties are still maintained after the application of the SSMG which ensures the derivation of *all redundant* association rules *without loss of information*. Unfortunately, this is not the case for the informative generic basis \mathcal{IGB} [7]. Indeed, even if it was proven in [7] that \mathcal{IGB} also verifies the ideal properties of an association rule representation, the obtained generic basis, once the SSMG is applied to \mathcal{IGB} , is with information loss because some *succinct* MGs can sometimes be missing (*w.r.t.* the definition of \mathcal{IGB} , see [7]).

The couple $(\mathcal{SGB}, \mathcal{SRI})$ of *succinct* generic bases of association rules is defined as follows ⁽⁴⁾:

Definition 12. (THE SUCCINCT GENERIC BASIS (\mathcal{SGB}) FOR EXACT ASSOCIATION RULES) *Let $\mathcal{FCI}_{\mathcal{K}}$ be the set of the frequent CIs extracted from a context \mathcal{K} . For each entry f in $\mathcal{FCI}_{\mathcal{K}}$, let $\mathcal{FMG}_{\text{SUC}f}$ be the set of its succinct frequent MGs. The succinct generic basis for exact association rules \mathcal{SGB} is given by: $\mathcal{SGB} = \{R: g \Rightarrow (f \setminus g) \mid f \in \mathcal{FCI}_{\mathcal{K}} \wedge g \in \mathcal{FMG}_{\text{SUC}f} \wedge g \neq f$ ⁽⁵⁾ $\}$.*

Definition 13. (THE SUCCINCT TRANSITIVE REDUCTION (\mathcal{SRI}) FOR APPROXIMATE ASSOCIATION RULES) *Let $\mathcal{FMG}_{\text{SUC}\mathcal{K}}$ be the set of the succinct frequent MGs extracted from a context \mathcal{K} . The succinct transitive reduction \mathcal{SRI} is given by: $\mathcal{SRI} = \{R: g \Rightarrow (f \setminus g) \mid f \in \mathcal{FCI}_{\mathcal{K}} \wedge g \in \mathcal{FMG}_{\text{SUC}\mathcal{K}} \wedge f \in \text{Cov}^u(f_1) \text{ with } f_1 = g'' \wedge \text{Conf}(R) = \frac{\text{Supp}(f)}{\text{Supp}(g)} \geq \text{minconf}\}$.*

Example 3. Consider the extraction context \mathcal{K} given by Figure 1 (Left) for a *minsupp* value equal to 1. The lexicographic order relation is used as a total one. The associated Iceberg concept lattice is depicted by Figure 2 (Left). A *succinct* exact generic rule is an “intra-node” association, with a confidence value equal to 1, within an equivalence class of the Iceberg concept lattice. The use of the SSMG allows, for example, to only

³ \mathcal{DGB} (resp. \mathcal{LB} , \mathcal{RR} , and \mathcal{NRR}) stands for Duquenne-Guigues Basis [4] (resp. Luxenburger Basis [4], Representative Rules [5], and Non-Redundant Rules [6]).

⁴ The definition of the couple $(\mathcal{GB}, \mathcal{RI})$ can be derived from that of the couple $(\mathcal{SGB}, \mathcal{SRI})$ by considering *all* MGs instead of only *succinct* ones.

⁵ The condition $g \neq f$ ensures discarding non-informative association rules of the form $g \Rightarrow \emptyset$.

extract the *succinct* exact generic association rule $adg \Rightarrow bce$ from the equivalence class having “abcdeg” for *frequent* CI, instead of four if *redundant frequent* MGs were of use (as indicated by the last entry in the table of Figure 1 (Right)). A *succinct* approximate generic rule represents an “inter-node” association, assorted with a confidence measure, between an equivalence class and another belonging to its upper cover. For example, for $minconf = 0.4$, only the association rule $ad \xrightarrow{0.5} bceg$ is extracted from both equivalence classes having, respectively, “abcde” and “abcdefg” for *frequent* CI instead of four if *redundant frequent* MGs were of use (as indicated by the seventh entry in the table of Figure 1 (Right)). The complete set of *succinct* generic association rules, extracted from \mathcal{K} , is reported in Figure 2 (Right). The cardinality of \mathcal{SGB} (resp. \mathcal{GB}) is equal to **13** (resp. **23**), while that of \mathcal{SRI} (resp. \mathcal{RI}) is equal to **21** (resp. **28**). Hence, thanks to the SSMG, we are able to discard **43.48%** (resp. **25%**) of the exact (resp. approximate) generic association rules since they are *redundant*. Note that the total number of association rules, which can be retrieved from \mathcal{K} , is equal to **943**.

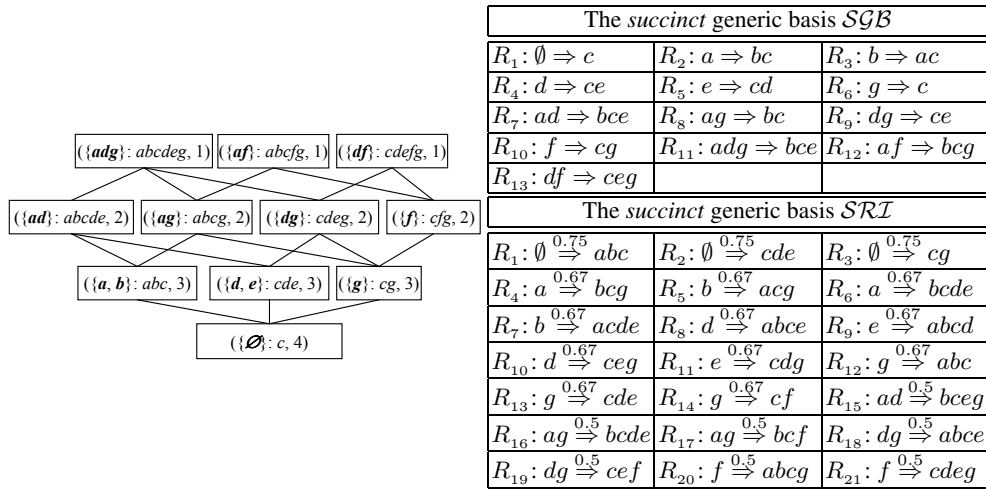


Fig. 2. (Left) For $minsupp = 1$, the Iceberg concept lattice associated to the extraction context \mathcal{K} of Figure 1 (Left). Each one of its equivalence classes contains a *frequent* CI f accompanied by the set of its *succinct frequent* MGs $FMG_{supp} f$ and its support, in the form $(FMG_{supp} f: f, Supp(f))$. (Right) The complete set of *succinct* generic association rules extracted from \mathcal{K} .

4.3 Derivation of redundant association rules

In the following, we study the structural properties of the new generic bases introduced in the previous subsection. The study requires checking the *ideal* properties of an association rule representation (see Definition 11). Since, it was shown in [5] that the couple $(\mathcal{GB}, \mathcal{RI})$ is extracted without loss of information, it is sufficient to show that it is possible to derive without loss of information *all* association rules that belong to the couple

$(\mathcal{GB}, \mathcal{RT})$ starting from the couple $(S\mathcal{GB}, S\mathcal{RT})$. If so, *all redundant* association rules can be derived from $(S\mathcal{GB}, S\mathcal{RT})$.

Association rules belonging to the couple $(S\mathcal{GB}, S\mathcal{RT})$ are implications between *succinct frequent* minimal generators (MGs) and *frequent* closed itemsets (CIs). Hence, to derive the couple $(\mathcal{GB}, \mathcal{RT})$, *redundant frequent* MGs need to be deduced since they form the premises of *redundant* generic association rules, *i.e.*, association rules belonging to $(\mathcal{GB}, \mathcal{RT})$ and discarded from $(S\mathcal{GB}, S\mathcal{RT})$. In order to derive *all* association rules belonging to $(\mathcal{GB}, \mathcal{RT})$, we propose a new axiom called the *substitution axiom*. Thus, from each association rule $R: X \Rightarrow (Y \setminus X)$ of $(S\mathcal{GB}, S\mathcal{RT})$ where $X \in \mathcal{FMG}_{\text{succ}\mathcal{K}}$ and $Y \in \mathcal{FCI}_{\mathcal{K}}$, we propose to derive, using the substitution axiom, the set of *redundant* generic association rules given by: $\text{Red_Gen_Assoc_Rules}_R: X \Rightarrow (Y \setminus X) = \{R': Z \Rightarrow (Y \setminus Z) \mid Z \in \mathcal{FMG}_{\text{red}\mathcal{K}} \text{ s.t. } X \models Z\}$. The substitution axiom proceeds according to the following steps:

Step 1 The set \mathcal{GB} (resp. \mathcal{RT}) is firstly initialized to $S\mathcal{GB}$ (resp. $S\mathcal{RT}$).

Step 2 The association rules belonging to $(\mathcal{GB}, \mathcal{RT})$ are processed in an ascending order of their respective sizes⁽⁶⁾, *i.e.*, that for an association rule $R: X \Rightarrow (Y \setminus X) \in (\mathcal{GB}, \mathcal{RT})$ where $X \in \mathcal{FMG}_{\text{succ}\mathcal{K}}$ and $Y \in \mathcal{FCI}_{\mathcal{K}}$, the set of the *redundant* generic association rules associated to each association rule $R_1: X_1 \Rightarrow (Y_1 \setminus X_1)$, such that $X_1 \subset X$ and $Y_1 \subset Y$, were already derived.

Step 2.1 For each association rule $R: X \Rightarrow (Y \setminus X) \in \mathcal{GB}$, derive the set of *redundant* generic association rules $\text{Red_Gen_Assoc_Rules}_R = \{R': Z \Rightarrow (Y \setminus Z) \mid Z \text{ is the result of the substitution of a subset of } X, \text{ say } V, \text{ by } T \text{ s.t. } (R_1: V \Rightarrow (I \setminus V), R_2: T \Rightarrow (I \setminus T)) \in \mathcal{GB} \text{ with } I \in \mathcal{FCI}_{\mathcal{K}} \text{ and } \nexists Z_1 \subseteq Z \text{ s.t. } Z_1 \Rightarrow (Y \setminus Z_1) \in \mathcal{GB}\}$.

Step 2.2 For each association rule $R: X \Rightarrow (Y \setminus X) \in \mathcal{RT}$, derive the set of *redundant* generic association rules $\text{Red_Gen_Assoc_Rules}_R = \{R': Z \Rightarrow (Y \setminus Z) \mid Z \text{ is the result of the substitution of a subset of } X, \text{ say } V, \text{ by } T \text{ s.t. } (R_1: V \Rightarrow (I \setminus V), R_2: T \Rightarrow (I \setminus T)) \in \mathcal{GB} \text{ with } I \in \mathcal{FCI}_{\mathcal{K}} \text{ and } \nexists Z_1 \subseteq Z \text{ s.t. } Z_1 \Rightarrow (Y \setminus Z_1) \in \mathcal{RT}\}$. ♦

Note that comparing Z to Z_1 ensures discarding the case where a substitution leads to an already existing association rule or to a one having a *non-minimal* generator as a premise.

Example 4. From the association rule $R: adg \Rightarrow bce$ belonging to $S\mathcal{GB}$ (*c.f.*, Figure 2 (Right)), we will show how to derive association rules belonging to \mathcal{GB} which are *redundant w.r.t.* R . Before that R is processed, *all* association rules whose respective sizes are lower than that of R (*i.e.*, lower than 6) were handled and *redundant* generic association rules were derived from such association rules. Among the handled association rules, we find those having for premises the 2-subsets of “ adg ”, *i.e.*, $ad \Rightarrow bce$, $ag \Rightarrow bc$ and $dg \Rightarrow ce$. To derive the *redundant* generic association rules associated to R , the first 2-subset of “ adg ”, *i.e.*, “ ad ”, is replaced by the *frequent* MGs having its closure, *i.e.*, the *redundant frequent* MGs “ ae ”, “ bd ” and “ be ”. Indeed, generic association rules using these latter as premises were already derived as *redundant w.r.t.* ad

⁶ The size of an association rule is equal to the number of items it contains.

$\Rightarrow bce$. Hence, we augment \mathcal{GB} by the following association rules: $aeg \Rightarrow bcd$, $bdg \Rightarrow ace$ and $beg \Rightarrow acd$. The same process is applied to the second subset of “ adg ”, i.e., “ ag ”. Nevertheless, the obtained association rule, namely $bdg \Rightarrow ace$, will not be added to \mathcal{GB} . Indeed, it already exists an association rule in \mathcal{GB} such that $Z_1 \Rightarrow (abcdeg \setminus Z_1)$ and $Z_1 \subseteq abg$ (Z_1 being itself equal to “ abg ”). It is the same for the derived association rule using the third subset “ dg ”, i.e., $aeg \Rightarrow bcd$ (Z_1 being equal to “ aeg ”).

Now, we prove that the substitution axiom allows the couple $(S\mathcal{GB}, S\mathcal{RI})$ to be *lossless* and *sound*. Then, we show that this couple is also *informative*.

Proposition 2. *The couple $(S\mathcal{GB}, S\mathcal{RI})$ of generic bases is lossless: $\forall R: X \Rightarrow (Y \setminus X) \in (S\mathcal{GB}, S\mathcal{RI})$, the set $\text{Red_Gen_Assoc_Rules}_R = \{R': Z \Rightarrow (Y \setminus Z) \mid Z \in \mathcal{FMG}_{\text{red}_K} \text{ s.t. } X \models Z\}$ of the redundant generic association rules with respect to R , is completely derived thanks to the proposed substitution axiom.*

Proof.

The sorting imposed in Step 2 ensures that, before R is processed, all association rules whose respective sizes are lower than that of R were handled, and redundant generic association rules were then derived from such association rules. Hence, all information required to derive association rules belonging to $\text{Red_Gen_Assoc_Rules}_R$ are gathered thanks to the different sets $\text{Red_Gen_Assoc_Rules}_{R_1}: X_1 \Rightarrow (Y_1 \setminus X_1)$ such that $X_1 \in \mathcal{FMG}_{\text{sup}_K}$, $Y_1 \in \mathcal{FCI}_K$ and $Y_1 \subset Y$. Indeed, using these sets, all redundant frequent MGs, with respect to X , are straightforwardly derivable since, for each subset X_1 of X , the different frequent MGs belonging to its equivalence class are already known as they are the premises of association rules belonging to the sets $\text{Red_Gen_Assoc_Rules}_{R_1}$ defined above (see Definition 7 and Definition 9 for the details on derivation). Hence, all association rules belonging to $(\mathcal{GB}, \mathcal{RI})$ can be deduced from $(S\mathcal{GB}, S\mathcal{RI})$ using the substitution axiom. Therefore, the couple $(S\mathcal{GB}, S\mathcal{RI})$ is *lossless*. \blacklozenge

Proposition 3. *The couple $(S\mathcal{GB}, S\mathcal{RI})$ of generic bases is sound: $\forall R': Z \Rightarrow (Y \setminus Z) \in \text{Red_Gen_Assoc_Rules}_R: X \Rightarrow (Y \setminus X)$, $\text{Supp}(R') = \text{Supp}(R)$ and $\text{Conf}(R') = \text{Conf}(R)$.*

Proof.

On the one hand, $\text{Supp}(R)$ is equal to $\text{Supp}(Y)$. It is the same for $\text{Supp}(R')$. Hence, $\text{Supp}(R') = \text{Supp}(R)$. On the other hand, X and Z are two frequent MGs belonging to the same equivalence class. Hence, $\text{Supp}(X)$ is equal to $\text{Supp}(Z)$. Thus, $\text{Conf}(R') = \frac{\text{Supp}(Y)}{\text{Supp}(Z)} = \frac{\text{Supp}(Y)}{\text{Supp}(X)} = \text{Conf}(R)$. Therefore, the couple $(S\mathcal{GB}, S\mathcal{RI})$ is *sound*. \blacklozenge

The property of derivability is verified by the couple $(S\mathcal{GB}, S\mathcal{RI})$ of generic bases since it is *lossless* and *sound*. Now, we show that this couple allows the retrieval of the exact values of the support and the confidence associated to each derived association rule.

Proposition 4. *The couple $(S\mathcal{GB}, S\mathcal{RI})$ of generic bases is informative: the support and the confidence of all derived association rules can exactly be retrieved from $(S\mathcal{GB}, S\mathcal{RI})$.*

Proof.

Association rules belonging to the couple $(S\mathcal{G}\mathcal{B}, S\mathcal{R}\mathcal{I})$ are of the following form: $g \Rightarrow (f \setminus g)$ where $g \in \mathcal{FMG}_{SUC\mathcal{K}}$ and $f \in \mathcal{FCL}_{\mathcal{K}}$. Therefore, we are able to reconstitute all necessary frequent CIs by concatenation of the premise and the conclusion parts of the generic association rules belonging to $(S\mathcal{G}\mathcal{B}, S\mathcal{R}\mathcal{I})$. Since the support of a frequent itemset I is equal to the support of the smallest frequent CI containing it [8], then the support of I and its closure can be straightforwardly derived from $(S\mathcal{G}\mathcal{B}, S\mathcal{R}\mathcal{I})$. Hence, the support and the confidence values of all redundant association rules can exactly be retrieved. Thus, the couple $(S\mathcal{G}\mathcal{B}, S\mathcal{R}\mathcal{I})$ is informative. \blacklozenge

The substitution axiom is proved to be lossless, sound and informative; allowing to derive *all* association rules forming $(\mathcal{G}\mathcal{B}, \mathcal{R}\mathcal{I})$ as well as their *exact* support and confidence values. Since the couple $(\mathcal{G}\mathcal{B}, \mathcal{R}\mathcal{I})$ is shown to be extracted without loss of information [5], we can deduce that the couple $(S\mathcal{G}\mathcal{B}, S\mathcal{R}\mathcal{I})$ is also extracted without information loss. In order to find the complete set of *valid redundant* association rules that can be extracted from a context \mathcal{K} , the axiom of transitivity proposed by Luxemburger [20] should be applied to the generic basis $\mathcal{R}\mathcal{I}$ to derive association rules forming the informative basis $\mathcal{I}\mathcal{B}$ for the approximate association rules [3]. Then, the cover operator proposed by Kryszkiewicz [5] or the lossless and sound axiomatic system proposed by Ben Yahia and Mephu Nguifo [21] makes it possible to derive *all valid redundant* association rules starting from the couple $(\mathcal{G}\mathcal{B}, \mathcal{I}\mathcal{B})$ of generic bases. The complete process allowing to derive *all valid (redundant)* association rules (denoted $\mathcal{A}\mathcal{R}$), starting from the couple $(S\mathcal{G}\mathcal{B}, S\mathcal{R}\mathcal{I})$, is hence as follows (the words *axiom* and *operator* are omitted):

$$(S\mathcal{G}\mathcal{B}, S\mathcal{R}\mathcal{I}) \xrightarrow{\text{substitution}} (\mathcal{G}\mathcal{B}, \mathcal{R}\mathcal{I}) \xrightarrow{\text{transitivity}} (\mathcal{G}\mathcal{B}, \mathcal{I}\mathcal{B}) \xrightarrow{\text{cover or Ben Yahia et al.}} \mathcal{A}\mathcal{R}.$$

5 Experimental study

We carried out experimentations on benchmark datasets⁽⁷⁾ in order to evaluate the number of (*succinct*) generic association rules. Characteristics of these datasets are summarized by Table 1. All experiments were run on a PC equipped with a 2.4GHz Pentium IV and 512MB of RAM. All programs were implemented in the C language and compiled with gcc 3.3.1 under the distribution S.U.S.E Linux 9.0. Hereafter, we use a logarithmically scaled ordinate axis in all figures.

We compared both couples $(S\mathcal{G}\mathcal{B}, S\mathcal{R}\mathcal{I})$ and $(\mathcal{G}\mathcal{B}, \mathcal{R}\mathcal{I})$ using the couple size as evaluation criterion, for a fixed *minsupp* value. Indeed, this was carried out for the PUMSB (resp. CONNECT, MUSHROOM and T40I10D100K) dataset for a *minsupp* value equal to **70%** (resp. **50%**, **0.01%** and **1%**). Obtained results are graphically sketched by Figure 3. For each dataset, the *minconf* value varies between the aforementioned *minsupp* value and **100%**.

Figure 3 points out that removing redundancy within the *frequent* MG set⁽⁸⁾ offers an interesting lossless reduction of the number of the extracted generic association rules.

⁷ These benchmark datasets are downloadable from: <http://fimi.cs.helsinki.fi/data>.

⁸ Interested readers are referred to [1] for more details.

Dataset	Number of items	Number of objects	Average object size	<i>minsupp</i> interval (%)
PUMSB	7, 117	49, 046	74	90 - 60
MUSHROOM	119	8, 124	23	1 - 0.01
CONNECT	129	67, 557	43	90 - 50
T40I10D100K	1, 000	100, 000	40	10 - 1

Table 1. Dataset characteristics.

Indeed, the use of the SSMG allows to remove in average **63.03%** (resp. **49.46%**) of the *redundant* generic association rules extracted from the PUMSB (resp. MUSHROOM) dataset. The maximum rate of redundancy reaches **68.11%** (resp. **53.84%**) for the PUMSB (resp. MUSHROOM) dataset, for a *minconf* value equal to **100%** (resp. **20%**). For the CONNECT and T40I10D100K datasets, the respective curves representing the size of the couple ($SG\mathcal{B}$, SRI) and those representing the size of the couple ($\mathcal{G}\mathcal{B}$, $\mathcal{R}\mathcal{I}$) are strictly overlapping. Indeed, these two datasets do not generate *redundant frequent* MGs and, hence, there are no *redundant* generic association rules. Furthermore, for the T40I10D100K dataset, none *exact* association rule is generated since each *frequent* MG is equal to its closure.

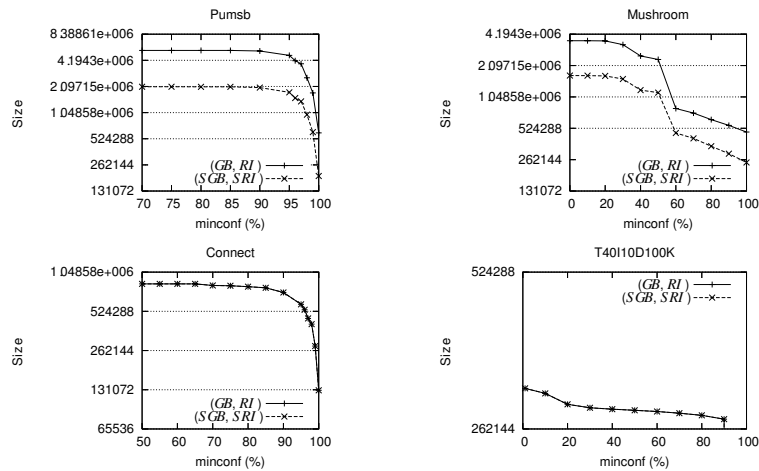


Fig. 3. For a fixed *minsupp* value, the size of the couple ($\mathcal{G}\mathcal{B}$, $\mathcal{R}\mathcal{I}$) of generic bases compared to that of the couple ($SG\mathcal{B}$, SRI) of *succinct* generic bases.

We also set the *minconf* value to **0%** to evaluate the reduction rate within *exact* generic association rules (*i.e.*, the generic basis $\mathcal{G}\mathcal{B}$) compared to that within *approximate* ones (*i.e.*, the generic basis $\mathcal{R}\mathcal{I}$). In this context, Figure 4 shows that, for the PUMSB dataset, in average **62.46%** (resp. **49.11%**) of the exact (resp. approximate)

generic association rules are *redundant*, and the maximum rate of redundancy reaches **68.46%** (resp. **62.65%**) for a *minsupp* value equal to **65%** (resp. **65%**). For the MUSHROOM dataset, in average **50.55%** (resp. **52.65%**) of the exact (resp. approximate) generic association rules are *redundant*, and the maximum rate of redundancy reaches **53.23%** (resp. **57.86%**) for a *minsupp* value equal to **0.20%** (resp. **0.10%**).

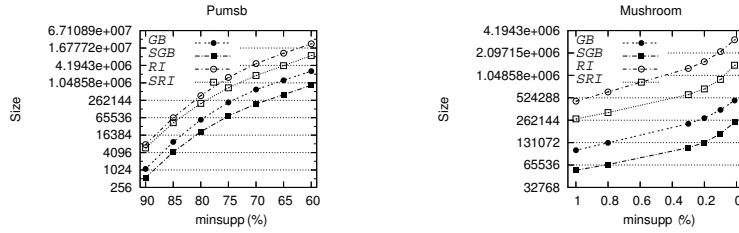


Fig. 4. For a fixed *minconf* value, the size of the generic basis \mathcal{GB} (resp. \mathcal{RI}) compared to that of the *succinct* generic basis \mathcal{SGB} (resp. \mathcal{SRI}).

These experiments clearly indicate that our approach can be advantageously used to eliminate, without loss of information, a large number of *redundant* generic association rules.

6 Conclusion and future work

In this paper, we briefly described the principal structural properties of the succinct system of minimal generators (SSMG) redefined in [1]. We then incorporated it into the framework of generic bases to tackle the problem of succinctness within generic association rule sets. Thus, we introduced two new *succinct* generic bases of association rules, namely the couple $(\mathcal{SGB}, \mathcal{SRI})$. We also showed that, starting from this couple, it is possible to derive without loss of information *all valid* association rules belonging to the couple $(\mathcal{GB}, \mathcal{RI})$ thanks to the application of a new substitution process. Consequently, any *valid redundant* association rule, which can be extracted from a context, can be inferred starting from the couple $(\mathcal{SGB}, \mathcal{SRI})$. Finally, carried out experiments confirmed that the application of the SSMG makes it possible to eliminate, as much as possible, *redundant* generic association rules and, hence, to only offer succinct and informative ones to users.

In the near future, we plan to set up an association rule visualization platform based on *succinct* generic bases, which, in our opinion, will constitute a helpful tool for the users. In this context, integrating the quality measures and the user-defined constraints in this tool will be interesting for further association rule pruning. In addition, we think that a careful study of the effect of the total order relation choice, on the quality of the extracted *succinct* association rules according to the data under consideration, presents an interesting issue towards increasing the knowledge usefulness.

Acknowledgements

We are grateful to the anonymous reviewers for their helpful comments. This work is partially supported by the French-Tunisian project CMCU 05G1412.

References

1. Hamrouni, T., Ben Yahia, S., Mephu Nguifo, E.: Succinct system of minimal generators: A thorough study, limitations and new definitions. In: Proceedings of the 4th International Conference on Concept Lattices and their Applications (CLA 2006), Hammamet, Tunisia. (2006)
2. Agrawal, R., Imielinski, T., Swami, A.: Mining association rules between sets of items in large databases. In: Proceedings of the ACM-SIGMOD International Conference on Management of Data, Washington D. C., USA. (1993) 207–216
3. Bastide, Y., Pasquier, N., Taouil, R., Stumme, G., Lakhal, L.: Mining minimal non-redundant association rules using frequent closed itemsets. In: Proceedings of the 1st International Conference on Computational Logic (DOOD 2000), Springer-Verlag, LNAI, volume 1861, London, UK. (2000) 972–986
4. Stumme, G., Taouil, R., Bastide, Y., Pasquier, N., Lakhal, L.: Intelligent structuring and reducing of association rules with formal concept analysis. In Baader, F., Brewker, G., Eiter, T., eds.: Proceedings of the Joint German/Austrian Conference on AI: Advances in Artificial Intelligence, Springer-Verlag, LNAI, volume 2174, Heidelberg, Germany. (2001) 335–350
5. Kryszkiewicz, M.: Concise representation of frequent patterns and association rules. Habilitation dissertation, Institute of Computer Science, Warsaw University of Technology, Warsaw, Poland. (2002)
6. Zaki, M.J.: Mining non-redundant association rules. *Journal of Data Mining and Knowledge Discovery (DMKD)*, volume 9. (2004) 223–248
7. Gasmi, G., BenYahia, S., Mephu Nguifo, E., Slimani, Y.: *TGB*: A new informative generic base of association rules. In: Proceedings of the International 9th Pacific-Asia Conference on Knowledge Data Discovery (PAKDD 2005), Springer-Verlag, LNAI, volume 3518, Hanoi, Vietnam. (2005) 81–90
8. Pasquier, N., Bastide, Y., Taouil, R., Lakhal, L.: Efficient mining of association rules using closed itemset lattices. *Journal of Information Systems*, volume 24 (1) (1999) 25–46
9. Calders, T., Goethals, B.: Mining all non-derivable frequent itemsets. In Elomaa, T., Mannila, H., Toivonen, H., eds.: Proceedings of the 6th European Conference on Principles and Practice of Knowledge Discovery in Databases (PKDD 2002), Springer-Verlag, LNCS, volume 2431, Helsinki, Finland. (2002) 74–85
10. Boulicaut, J.F., Bykowski, A., Rigotti, C.: Free-sets: A condensed representation of boolean data for the approximation of frequency queries. In *Journal of Data Mining and Knowledge Discovery (DMKD)*, volume 7(1). (2003) 5–22
11. Ohsaki, M., Kitaguchi, S., Okamoto, K., Yokoi, H., Yamaguchi, T.: Evaluation of rule interestingness measures with a clinical dataset on Hepatitis. In Boulicaut, J.F., Esposito, F., Giannotti, F., Pedreschi, D., eds.: Proceedings of the 8th European Conference on Principles and Practice of Knowledge Discovery in Databases (PKDD 2004), Springer-Verlag, LNCS, volume 3202, Pisa, Italy. (2004) 362–373
12. Srikant, R., Vu, Q., Agrawal, R.: Mining association rules with item constraints. In: Proceedings of the 3rd International Conference on Knowledge Discovery in Databases and Data Mining (KDD 1997), Newport Beach, California, USA. (1997) 67–73

13. Bonchi, F., Lucchese, C.: On condensed representations of constrained frequent patterns. In *Journal of Knowledge and Information Systems* (2005) 1–22
14. Wille, R.: Restructuring lattices theory: An approach based on hierarchies of concepts. I. Rival, editor, *Ordered Sets*, Reidel, Dordrecht-Boston (1982) 445–470
15. Kryszkiewicz, M.: Representative association rules and minimum condition maximum consequence association rules. In: *Proceedings of 2nd European Symposium on Principles of Data Mining and Knowledge Discovery (PKDD 1998)*, Springer-Verlag, LNCS, volume 1510, Nantes, France. (1998) 361–369
16. Dong, G., Jiang, C., Pei, J., Li, J., Wong, L.: Mining succinct systems of minimal generators of formal concepts. In: *Proceedings of the 10th International Conference on Database Systems for Advanced Applications (DASFAA 2005)*, Springer-Verlag, LNCS, volume 3453, Beijing, China. (2005) 175–187
17. Ganter, B., Wille, R.: *Formal Concept Analysis*. Springer-Verlag (1999)
18. Bastide, Y., Taouil, R., Pasquier, N., Stumme, G., Lakhal, L.: Mining frequent patterns with counting inference. In: *Proceeding of the 6th ACM-SIGKDD International Conference on Knowledge Discovery and Data Mining*, volume 2(2), Boston, Massachusetts, USA. (2000) 66–75
19. Agrawal, R., Srikant, R.: Fast algorithms for mining association rules. In Bocca, J.B., Jarke, M., Zaniolo, C., eds.: *Proceedings of the 20th International Conference on Very Large Databases*, Santiago, Chile. (1994) 478–499
20. Luxenburger, M.: Implication partielles dans un contexte. *Mathématiques et Sciences Humaines*, volume 29 (113) (1991) 35–55
21. Ben Yahia, S., Mephu Nguifo, E.: Revisiting generic bases of association rules. In: *Proceedings of 6th International Conference on Data Warehousing and Knowledge Discovery (DAWAK 2004)*, Springer-Verlag, LNCS, volume 3181, Zaragoza, Spain. (2004) 58–67