

GARC_M: Generic Association Rules based Classifier Multi-parameterizable

I. Bouzouita, S. Elloumi, and S. Ben Yahia

Faculty of Sciences of Tunis,
Computer Science Department, 1060 Tunis, Tunisia.
{samir.elloumi;sadok.benyahia}@fst.rnu.tn

Abstract. Many studies in data mining have proposed a new classification approach called *associative classification*. According to several reports associative classification achieves higher classification accuracy than do traditional classification approaches. However, the associative classification suffers from a major drawback: it is based on the use of a very large number of classification rules; and consequently takes efforts to select the best ones in order to construct the classifier. To overcome such drawback, we propose a new associative classification method called *GARC_M* that exploits a generic basis of association rules in order to reduce the number of association rules without jeopardizing the classification accuracy. Moreover, *GARC_M* proposes to users some interestingness measures that arise from data mining in order to select the best rules during classification of new instances. Carried out experiments on 12 benchmark data sets indicate that *GARC_M* is highly competitive in terms of accuracy in comparison with popular associative classification methods.

Keywords: Associative Classification, Generic Basis, Classification Rules, Generic association rules, Classifier.

1 Introduction

In the last decade, a new approach called *associative classification* (AC) was proposed to integrate association rule mining and classification in order to handle large databases. Given a training data set, the task of an associative classification algorithm is to discover the classification rules which satisfy the user specified constraints denoted respectively by minimum support (*minsup*) and minimum confidence (*minconf*) thresholds. The classifier is built by choosing a subset of the generated classification rules that could be of use to classify new objects or instances. Many studies have shown that AC often achieves better accuracy than do traditional classification techniques [1, 2]. In fact, it could discover interesting rules omitted by well known approaches such as C4.5 [3]. However, the main drawback of this approach is that the number of generated associative classification rules could be large and takes efforts to retrieve, prune, sort and select high quality rules among them. To overcome this problem, we propose a new approach called *GARC_M* which uses generic bases of association rules. The

main originality of $GARC_M$ is that it extracts the generic classification rules directly from a generic basis of association rules, in order to retain a small set of rules with higher quality and lower redundancy in comparison with current AC approaches. Moreover, a new score is defined by the $GARC_M$ approach to find an effective rule selection during the class label prediction of a new instance, in the sake of reducing the error rate. This tackled issue is quite challenging, since the goal is to use generic rules while maintaining a high classifier accuracy.

The remainder of the paper is organized as follows. Section 2 briefly reports basic concepts of associative classification and scrutinizes related pioneering works. Generic bases of association rules are surveyed in section 3. Section 4 presents our proposed approach, where details about classification rules discovery, building classifier and prediction of test instances are discussed. Experimental results and comparisons are given in section 5. Finally, section 6 concludes this paper and points out future perspectives.

2 Associative Classification

An association rule is a relation between itemsets having the following form: $R : X \Rightarrow Y - X$, where X and Y are frequent itemsets for a minimal support $minsup$, and $X \subset Y$. Itemsets X and $(Y - X)$ are called, respectively, *premise* and *conclusion* of the rule R . An association rule is valid whenever its strength metric, $confidence(R) = \frac{support(Y)}{support(X)}$, is greater than or equal to the minimal threshold of confidence $minconf$.

An associative classification rule (ACR) is a special case of an association rule. In fact, an ACR conclusion part is reduced to a single item referring a class attribute. For example, in an ACR such as $X \Rightarrow c_i$, c_i must be a class attribute.

2.1 Basic notions

Let us define the classification problem in an association rule task. Let D be the training set with n attributes (columns) A_1, \dots, A_n and $|D|$ rows. Let C be the list of class attributes.

Definition 1. *An object or instance in D can be described as a combination of attribute names and values a_i and an attribute class denoted by c_i [4].*

Definition 2. *An item is described as an attribute name and a value a_i [4].*

Definition 3. *An itemset can be described as a set of items contained in an object.*

A classifier is a set of rules of the form $A_1, A_2, \dots, A_n \Rightarrow c_i$ where A_i is an attribute and c_i is a class attribute. The classifier should be able to predict, as accurately as possible, the class of an unseen object belonging to the test data set. In fact, it should maximise the equality between the predicted class and the hidden actual class.

The AC achieves higher classification accuracy than do traditional classification approaches [1, 2]. The classification model is a set of rules easily understandable by humans and that can be edited [1, 2].

2.2 Related work

One of the first algorithms to use association rule approach for classification was CBA [4]. CBA, firstly, generates all the association rules with certain support and confidence thresholds as candidate rules by implementing the Apriori algorithm [5]. Then, it selects a small set from them by evaluating all the generated rules against the training data set. When predicting the class attribute for an example, the highest confidence rule, whose the body is satisfied by the example, is chosen for prediction.

CMAR [6] generates rules in a similar way as CBA with the exception that CMAR introduces a CR-tree structure to handle the set of generated rules and uses a set of them to make a prediction using a weighted χ^2 metric [6]. The latter metric evaluates the correlation between the rules.

ARC-AC and ARC-BC have been introduced in [7, 8] in the aim of text categorization. They generate rules similar to the Apriori algorithm and rank them in the same way as do CBA rules ranking method. ARC-AC and ARC-BC calculate the average confidence of each set of rules grouped by class attribute in the conclusion part and select the class attribute of the group with the highest confidence average.

The CPAR [2] algorithm adopts FOIL [9] strategy in generating rules from data sets. It seeks for the best rule itemset that brings the highest gain value among the available ones in data set. Once the itemset is identified, the examples satisfying it will be deleted until all the examples of the data set are covered. The searching process for the best rule itemset is a time consuming process, since the gain for every possible item needs to be calculated in order to determine the best item gain. During rule generation step, CPAR derives not only the best itemset but all close similar ones. It has been claimed that CPAR improves the classification accuracy whenever compared to popular associative methods like CBA and CMAR [2].

A new AC approach called Harmony was proposed in [10]. Harmony uses an instance-centric rule generation to discover the highest confidence discovering rules. Then, Harmony groups the set of rules into k groups according to their rule conclusions, where k is the total number of distinct class attributes in the training set. Within the same group of rules, Harmony sorts the rules in the same order as do CBA. To classify a new test instance, Harmony computes a score for each group of rules and assign the class attribute with the highest score or a set of class attributes if the underlying classification is a multi-class problem. It has been claimed that Harmony improves the efficiency of the rule generation process and the classification accuracy if compared to CPAR [2].

The main problem with AC approaches is that they generate an overwhelming number of rules during the learning stage. In order to overcome this drawback, our proposed approach tries to gouge this fact by the use of generic bases

of association rules in the classification framework. In the following, we begin by recall some key notions about the Formal Concept Analysis (FCA), a mathematical tool necessary for the derivation of generic bases of association rules.

3 GARC_M: Generic Association Rules based Classifier Multi-parameterizable

In this section, we propose a new AC method *GARC_M* that extracts the generic classification rules directly from a generic basis of association rules in order to overcome the drawback of the current AC approaches, i.e., the generation of a large number of associative classification rules. In the following, we will present the generic basis and then we will explain in details the *GARC_M* approach.

3.1 Generic Bases

The problem of the relevance and usefulness of extracted association rules is of primary importance. Indeed, in most real life databases, thousands and even millions of highly confident rules are generated among which many are redundant. In the following, we are interested in the lossless information reduction of association rules, which is based on the extraction of a generic subset of all association rules, called *generic basis* from which the remaining (redundant) association rules may be derived. In the following, we will present the generic basis of *IGB* [11].

The *IGB* basis is defined as follows:

Definition 4. Let $\mathcal{FCI}_{\mathcal{K}}$ be the set of frequent closed itemsets and \mathcal{G}_f be the set of minimal generators of all the frequent itemsets included or equal to a closed frequent itemset f . The *IGB* basis is defined as follows [11]:

$$\mathcal{IGB} = \{R: g_s \Rightarrow (f_1 - g_s) \mid f, f_1 \in \mathcal{FCI}_{\mathcal{K}} \text{ and } (f - g_s) \neq \emptyset \text{ and } g_s \in \mathcal{G}_f \wedge f_1 \preceq f \wedge \text{confidence}(R) \geq \text{minconf} \wedge \nexists g' \subset g_s \text{ such that } \text{confidence}(g' \Rightarrow f_1 - g') \geq \text{minconf}\}.$$

IGB has been shown to be informative and more compact than other generic basis.

3.2 Rule Generation

In this step, *GARC_M* extracts the generic basis of association rules. Once obtained, generic rules are filtered out to retain only rules whose conclusions include a class attribute. Then, by applying the decomposition axiom, we obtain new rules of the form $A_1, A_2, \dots, A_n \Rightarrow c_i$. Even though, the obtained rules are redundant, their generation is mandatory to guarantee a maximal cover of the necessary rules.

The *IGB* basis is composed of rules with a small premise which is an advantage for the classification framework when the rules imply the same class. For

example, let us consider two rules $R_1: A B C D \Rightarrow c11$ and $R_2: B C \Rightarrow c11$. R_1 and R_2 have the same attribute conclusion. R_2 is considered to be more interesting than R_1 , since it is needless to satisfy the properties $A D$ to choose the class $c11$. Hence, R_2 implies less constraints and can match more objects of a given population than R_1 .

Let us consider a new object $O_x: B C D$. If we have in the classifier just the rule R_1 , we cannot classify O_x because the attribute A does not permit the matching. However, the rule R_2 , which has a smaller premise than R_1 , can classify O_x . This example shows the importance of the generic rules and, especially, the use of the *TGB* basis to extract the generic classification rules. In fact, such set of rules is smaller than the number of all the classification rules and their use is beneficial for classifying new objects.

3.3 Classifier Builder

Unlike the current associative classification approaches, i.e., CBA, CMAR, ARC-AC and ARC-BC and Harmony, $GARC_M$ uses the generic classification rules to learn the classifier without setting any order on them. The major difference with current AC approaches [4, 6–8, 10] is that we use generic ACR directly deduced from generic bases of association rules to learn the classifier as shown by algorithm 1.

<pre> Data: \mathcal{D}: Training data, \mathcal{GR}: a set of generic classification rules Results: \mathcal{C}: Classifier Begin Foreach rule $r \in \mathcal{GR}$ do Foreach object $d \in \mathcal{D}$ do If d matches r.premise then remove d from \mathcal{D} and mark r if it correctly classifies d; If r is marked then insert r at the end of \mathcal{C}; select a default class; add the default class at the end of the classifier; return Classifier \mathcal{C} ; End </pre>

Algorithm 1: GARC: selected generic rules based on database coverage

3.4 New instance classification

After a set of rules is selected for classification, $GARC_M$ is ready to classify new objects. Some methods such as those described in [4, 7, 8, 10] are based on the support-confidence order to classify a new object. However, the confidence measure selection could be misleading, since it may identify a rule $A \Rightarrow B$ as an

interesting one even though, the occurrence of A does not imply the occurrence of B [12]. In fact, the confidence can be deceiving since it is only an estimate of the conditional probability of itemset B given an itemset A and does not measure the actual strength of the implication between A and B . Let us consider the example shown in Table 1 which shows the association between an item A and a class attribute B . A and \bar{A} represent respectively the presence and absence of item A , B represents a class attribute and \bar{B} the complement of B . We consider the associative classification $A \Rightarrow B$. The confidence of this rule is given by $confidence(A \Rightarrow B) = \frac{support(AB)}{support(A)} = \frac{201}{250} = 80.4\%$. Hence, this rule has high confidence.

In the following, we will introduce interestingness measures of association rules and give a semantic interpretation for each of them.

a. Lift or Interest The lift metric [12] computes the correlation between A and B as follows:

$$lift(A \Rightarrow B) = \frac{support(AB)}{support(A) * support(B)} = \frac{0.201}{0.250 * 0.900} = 0.893.$$

The fact that this quantity is less than 1 indicates negative correlation between A and B .

If the resulting value is greater than 1, then A and B are said positively correlated. If the resulting value is equal to 1, then A and B are independent and there is no correlation between them.

	B	\bar{B}	Total
A	201	49	250
\bar{A}	699	51	750
Total	900	100	1000

Table 1. Example

b. Least Confidence (or Surprise) The least confidence (or surprise) [13] metric is computed as follows:

$$Surprise(A \Rightarrow B) = (support(AB) - support(A\bar{B})) / support(B)$$

$$\text{logical rule: } surprise(A \Rightarrow B) = P(A) / P(B)$$

$$A \text{ and } B \text{ independent: } surprise(A \Rightarrow B) = 2 P(A) - (P(A) / P(B))$$

$$A \text{ and } B \text{ incompatible: } surprise(A \Rightarrow B) = - P(A) / P(B)$$

Surprise metric selects rules, even with small support value, having the premise A always with the conclusion B and nowhere else.

d. Loevinger Loevinger metric [13] is computed as follows:

$$loevinger(A \Rightarrow B) = (P(B/A) - P(B)) / P(\bar{B})$$

Unlike confidence metric, Loevinger metric does not suffer from the problem of producing misleading rules.

Based on the above study measures, we define a new lift based score formula as follows:

$$Score = \frac{1}{|Premise|} * lift^{\frac{|Premise|}{numberofitems}}$$

$$= \frac{1}{|Premise|} * \left(\frac{support(Rule)}{support(Premise) * support(Conclusion)} \right)^{\frac{|Premise|}{numberofitems}}$$

The introduced score includes the lift metric. In fact, the lift finds interesting relationships between A and B. It computes the correlation between the occurrence of A and B by measuring the real strength of the implication between them which is interesting for the classification framework. Moreover, the lift is divided by the cardinality of the rule premise part in order to give a preference to rules with small premises. Thus, $GARC_M$ collects the subset of rules matching the new object attributes from the classifier. Trivially, if all the rules matching it have the same class, $GARC_M$ just assigns that class to the new object. If the rules do not imply the same class attribute, the score firing is computed for each rule. The rule with the highest score value is selected to classify the new object.

4 Experiments

We have conducted experiments to evaluate the accuracy of our proposed approach $GARC_M$, developed in C++, and compared it to the well known classifiers CBA, ID3, C4.5 and Harmony. Experiments were conducted using 12 data sets taken from UCI Machine Learning Repository⁽¹⁾. The chosen data sets were discretized using the LUCS-KDD⁽²⁾ software.

The features of these data sets are summarized in Table 2. All the experiments were performed on a 2.4 GHz Pentium IV PC under Redhat Linux 7.2.

Data set	# attributes	# transactions	# classes
MONKS1	6	124	2
MONKS2	6	169	2
MONKS3	6	122	2
SPECT	23	80	2
PIMA	38	768	2
TIC TACTOE	29	958	2
ZOO	42	101	7
IRIS	19	150	3
WINE	68	178	3
GLASS	48	214	7
FLARE	39	1389	9
PAGEBLOCKS	46	5473	5

Table 2. data set description

¹ Available at <http://www.ics.uci.edu/~mllearn/MLRepository.html>

² Available at <http://www.csc.liv.ac.uk/~frans/KDD/Software/LUCS-KDD-DN/lucs-kdd DN.html>

Classification accuracy can be used to evaluate the performance of classification methods. It is the percentage of correctly classified examples in the test set and can be measured by splitting the data sets into a training set and a test set.

During experiments, we have used available test sets for data sets Monks1, Monks2 and Monks3 and we applied the 10 cross-validation for the rest of data sets, in which a data set is divided into 10 subsets; each subset is in turn used as testing data while the remaining data is used as the training data set; then the average accuracy across all 10 trials is reported.

The parameters are set as the following. In the rule generation algorithm, *minsup* is set to 10% and *minconf* to 80%. In order to extract generic association rules, we used the PRINCE algorithm [15] to generate *IGB* basis.

To evaluate C4.5 and ID3, we used the WEKA⁽³⁾ software and the Harmony prototype was kindly provided by its authors. We have implemented the CBA algorithm in C++ under Linux.

In the following, we will compare the effectiveness of using different interestingness measures of association rules for the classification framework. For this, we conducted experiments with reference to accuracy in order to compare the measures impact while classifying new instances.

4.1 Evaluating measures impact

Data set	Surprise	Loevinger	Lift	Score
MONKS1	42.6	62.5	59.2	92.0
MONKS2	67.1	59.0	49.3	56.0
MONKS3	97.2	92.8	56.7	96.3
SPECT	67.0	67.0	67.0	68.9
PIMA	72.9	72.9	72.9	73.0
TICTACTOE	63.0	63.0	63.0	73.0
ZOO	83.0	83.0	67.2	90.0
IRIS	94.0	89.3	95.3	95.4
WINE	92.8	81.1	88.3	89.8
GLASS	52.0	52.0	52.0	64.0
FLARE	84.7	84.6	84.7	85.0
PAGEBLOCKS	89.7	89.7	89.7	89.8

Table 3. Evaluating measures *vs* accuracy

Table 3 represents a comparison between the accuracy given by the measures used by *GARC_M* while classifying new instances.

Table 3 points out that the use of the score firing permits to achieve the best accuracy for eight data sets among eleven. The use of the surprise measure permits to achieve the best accuracy for three data sets. We can conclude that a

³ Available at <http://www.cs.waikato.ac.nz/ml/Weka>

multi-parameterizable tool will be efficient to present for users in order to choose the best measure suitable for the studied data set.

In the following, we put the focus on comparing $GARC_M$ accuracy by using the score firing versus that of the well known classifiers ID3, C4.5, CBA and Harmony.

4.2 Generic classification rules impact

Data set	ID3	C4.5	CBA	Harmony	$GARC_M$
MONKS1	77.0	75.0	92.0	83.0	92.0
MONKS2	64.0	65.0	56.0	48.0	56.0
MONKS3	94.0	97.0	96.3	82.0	96.3
SPECT	65.0	64.0	67.0	-	68.9
PIMA	71.3	72.9	73.0	73.0	73.0
TICTACTOE	83.5	85.6	63.1	81.0	65.0
ZOO	98.0	92.0	82.2	90.0	90.0
IRIS	94.0	94.0	95.3	94.7	95.4
WINE	84.8	87.0	89.5	63.0	89.8
GLASS	64.0	69.1	52.0	81.5	64.0
FLARE	80.1	84.7	85.0	83.0	85.0
PAGEBLOCKS	92.3	92.4	89.0	91.0	89.8

Table 4. Accuracy comparison of ID3, C4.5, CBA, Harmony and $GARC_M$ algorithms

Table 4 represents the accuracy of the classification systems generated by ID3, C4.5, CBA, Harmony and $GARC_M$ on the twelve benchmark data sets. The best accuracy values obtained for each of data sets is highlighted in bold print. Table 4 shows that $GARC_M$ outperforms the traditional classification approaches, *i.e.*, ID3 and C4.5 on six data sets and the associative classification approaches on nine data sets.

Statistics depicted by Table 4 confirm the fruitful impact of the use of the generic rules. The main reason for this is that $GARC_M$ classifier contains generic rules with small premises. In fact, this kind of rule allows to classify more objects than those with large premises.

5 Conclusion

In this paper, we introduced a new classification approach called $GARC_M$ that aims to prune the set of classification rules without jeopardizing the accuracy and even ameliorates the predictive power by investigating interestingness measures. To this end, $GARC_M$ uses generic bases of association rules to drastically reduce the number of associative classification rules. Moreover, it proposes a new score to ameliorate the rules selection for unseen objects. Carried out experiments

outlined that $GARC_M$ is highly competitive in terms of accuracy in comparison with popular classification methods. In the near future, we will investigate new metrics for the rule selection and we will apply $GARC_M$ approach to a wide range of applications like text categorization and biological applications.

Acknowledgements We are deeply grateful to Frans Coenen at the university of Liverpool for providing us the discretized UCI data sets and addressing our questions. We also thank Jianyong Wang for providing us the Harmony executable code.

References

1. Zaiane, O., Antonie, M.: On pruning and tuning rules for associative classifiers. In: Ninth International Conference on Knowledge Based Intelligence Information And Engineering Systems (KES'05), Melbourne, Australia (2005) 966–973
2. Xiaoxin Yin, J.H.: CPAR: Classification based on Predictive Association Rules. In: Proceedings of the SDM, San Francisco, CA (2003) 369–376
3. Quinlan, J.R.: C4.5 : Programs for Machine Learning. (1993)
4. Liu, B., Hsu, W., Ma, Y.: Integrating classification and association rule mining. In: Knowledge Discovery and Data Mining. (1998) 80–86
5. Agrawal, R., Srikant, R.: Fast algorithms for mining association rules. In Bocca, J.B., Jarke, M., Zaniolo, C., eds.: Proceedings of the 20th Intl. Conference on Very Large Databases, Santiago, Chile. (1994) 478–499
6. Li, W., Han, J., Pei, J.: CMAR: Accurate and efficient classification based on multiple class-association rules. In: Proceedings of IEEE International Conference on Data Mining (ICDM'01), San Jose, CA, IEEE Computer Society (2001) 369–376
7. Antonie, M., Zaiane, O.: Text Document Categorization by Term Association. In: Proc. of the IEEE International Conference on Data Mining (ICDM'2002), Maebashi City, Japan (2002) 19–26
8. Antonie, M., Zaiane, O.: Classifying Text Documents by Associating Terms with Text Categories . In: Proc. of the Thirteenth Austral-Asian Database Conference (ADC'02), Melbourne, Australia (2002)
9. Quinlan, J., Cameron-Jones, R.: FOIL: A midterm report. In: Proceedings of European Conference on Machine Learning, Vienna, Austria. (1993) 3–20
10. Wang, J., Karypis, G.: HARMONY: Efficiently mining the best rules for classification. In: Proceedings of the International Conference of Data Mining (SDM'05). (2005)
11. Gasmi, G., BenYahia, S., Nguifo, E.M., Slimani, Y.: $IG\mathcal{B}$: A new informative generic base of association rules. In: Proceedings of the Intl. Ninth Pacific-Asia Conference on Knowledge Data Discovery (PAKDD'05), LNAI 3518, Hanoi, Vietnam, Springer-Verlag (2005) 81–90
12. Han, J., Kamber, M.: Data Mining : Concepts and Techniques. Morgan Kaufmann. (2001)
13. Lallich, S., Teytaud, O.: Evaluation et validation de lintrt des rgles d'association. In: RNTI-E. (2004) 193–217
14. Totohasina, A., Ralambondrainy, H., Diatta, J.: Un algorithme efficace d'extraction des rgles d'association implicative. In: Proceedings of IEEE International Conference on Data Mining (ICDM'01), Hammamet, Tunisie (2004)

15. Hamrouni, T., BenYahia, S., Slimani, Y.: PRINCE : An algorithm for generating rule bases without closure computations. In Tjoa, A.M., Trujillo, J., eds.: Proceedings of 7th International Conference on Data Warehousing and Knowledge Discovery (DaWaK 2005), Springer-Verlag, LNCS 3589, Copenhagen, Denmark. (2005) 346–355