

# A qualitative methodology to reduce features in classification problems

Alvarez, M.A.<sup>a</sup>; Gonzalez-Abril, L.<sup>b</sup>; Ortega, J.A.<sup>a</sup>; Soria, L.M.<sup>a</sup>

<sup>a</sup>Department of Computer Science, University of Seville

<sup>b</sup>Department of Applied Economics I, University of Seville  
Seville, Spain

{maalvarez, luisgon, jortega, lsoria}@us.es

## Abstract

In this paper, a preliminary methodology which quantifies the dependence between features in a data set by using the Ameva discretization algorithm and the advantages of a qualitative model is developed. Thus, different matrices of interdependence are built providing a grade of dependence between two features. This methodology is applied to a well-known data set, obtaining promising results for the carried out system.

## 1 Introduction

The problem of classification is one of the main problems in data analysis and pattern recognition that requires the construction of a classifier, that is, a function that assigns a class label to instances described by a set of features. The induction of classifiers from data sets of classified instances is a central problem in machine learning. For that purpose, a large number of methodologies based on SVM [1], Naive Bayesian [2], C5.0 [3], etc. have been developed.

Additionally, qualitative modeling and reasoning is a very interesting area for applying and experimenting with machine learning techniques. Qualitative reasoning has special interest to systems where machine learning can be applied as modeling, diagnosis, control, discovery, design, and knowledge compilation.

One of the most important preprocess in classification is the discretization. This process establishes a relationship between continuous variables and their discrete transformation through functions. Therefore, it is possible to model qualitatively a series of continuous values if a label is assigned to them. Some studies [4] have shown that execute a prior process to discretize continuous features is more efficient than work directly with the continuous values. This process reduces the computation time and memory usage in the application of classification algorithms and it is used to manage the set of values of a feature more effectively. Some relevant discretization methods are Ameva [5], Khiops [6], CAIM [7] and others [8; 9].

The Ameva discretization method has been confirmed as one of the most promising algorithms due to its reduced execution time and the smaller number of intervals provided. This behavior is outstanding when the data set has a large

number of classes, although it has a slight reduction in the capacity of identification [5; 10].

Another problem in the classification process is the existence of irrelevant features [11]. When data is obtained experimentally, is not considered what features are relevant for the studied system. Several techniques [12; 13; 14] have been developed to reduce the number of features and to determine which are relevant for the system. Some of these techniques are based on principals components analysis [15] or factorial analysis [16].

The Ameva discretization algorithm [10] performs the discretization process effectively and quickly, so the set of values of a feature is greatly reduced, but do not reduce the number of features. Because Ameva uses the statistic  $\chi^2$  to determine the relationship between features and classes, it is possible to use this algorithm to determine the relationship between features.

In this paper, a new methodology based on Ameva algorithm is developed in order to reduce the number of features of a data set. This method exploits the advantages of Ameva in runtime and brings a different approach which was developed on.

The rest of this paper is organized as follows: first, the definition of the problem is presented in Section 2 to establish the notation of the rest of the paper. Also, the Ameva discretization algorithm and the Entropy coefficient are presented. Section 3 presents the new methodology to determine the dependence between features using the Ameva algorithm and the entropy coefficient. Section 4 reports the obtained results of applying the methodology in a toy example. The paper is finally concluded with a summary of the most important points and future works.

## 2 Discretization

Let  $X = \{x_1, x_2, \dots, x_N\}$  be a data set of a continuous attribute  $\mathcal{X}$  of mixed-mode data such that each example  $x_i$  belongs to only one of  $\ell$  classes of the variable denoted by

$$\mathcal{C} = \{C_1, C_2, \dots, C_\ell\}, \quad \ell \geq 2 \quad (1)$$

A continuous attribute discretization is a function  $\mathcal{D} : \mathcal{X} \rightarrow \mathcal{C}$  which assigns a class  $C_i \in \mathcal{C}$  to each value  $x \in \mathcal{X}$  in the domain of the property that is being discretized.

Let us consider a discretization  $\mathcal{D}$  which discretizes  $\mathcal{X}$  into  $k$  discrete intervals:

$$\mathcal{L}(k; X; \mathcal{C}) = \{L_1, L_2, \dots, L_k\}$$

where  $L_1$  is the interval  $[d_0, d_1]$  and  $L_j$  is the interval  $(d_{j-1}, d_j]$ ,  $j = 2, 3, \dots, k$ . Thus, a discretization variable is defined as  $\mathcal{L}(k) = \mathcal{L}(k; X; \mathcal{C})$  which verifies that, for all  $x_i \in X$ , a unique  $L_j$  exists such that  $x_i \in L_j$  for  $i = 1, 2, \dots, N$  and  $j = 1, 2, \dots, k$ . The discretization variable  $\mathcal{L}(k)$  of attribute  $\mathcal{X}$  and the class variable  $\mathcal{C}$  are treated from a descriptive point of view. Having two discrete attributes, a two-dimensional frequency table (called contingency table) as shown in the Table 1 can be built.

$C_i   L_j$	$L_1$	$\dots$	$L_j$	$\dots$	$L_k$	$n_{i.}$
$C_1$	$n_{11}$	$\dots$	$n_{1j}$	$\dots$	$n_{1k}$	$n_{1.}$
$\vdots$	$\vdots$	$\ddots$	$\vdots$	$\ddots$	$\vdots$	$\vdots$
$C_i$	$n_{i1}$	$\dots$	$n_{ij}$	$\dots$	$n_{ik}$	$n_{i.}$
$\vdots$	$\vdots$	$\ddots$	$\vdots$	$\ddots$	$\vdots$	$\vdots$
$C_\ell$	$n_{\ell 1}$	$\dots$	$n_{\ell j}$	$\dots$	$n_{\ell k}$	$n_{\ell.}$
$n_{.j}$	$n_{.1}$	$\dots$	$n_{.j}$	$\dots$	$n_{.k}$	$N$

Table 1: Contingency table

In Table 1,  $n_{ij}$  denotes the total number of continuous values belonging to the  $C_i$  class that are within the interval  $L_j$ .  $n_{i.}$  is the total number of instances belonging to the class  $C_i$ , and  $n_{.j}$  is the total number of instances that belong to the interval  $L_j$ , for  $i = 1, 2, \dots, \ell$  and  $j = 1, 2, \dots, k$ . So that:

$$n_{i.} = \sum_{j=1}^k n_{ij}, \quad n_{.j} = \sum_{i=1}^{\ell} n_{ij}, \quad N = \sum_{i=1}^{\ell} \sum_{j=1}^k n_{ij}$$

## 2.1 The Ameva discretization

Given discrete attributes  $\mathcal{C}$  and  $\mathcal{L}(k)$ , the contingency coefficient, denoted by  $\chi^2(k) \stackrel{def}{=} \chi^2(\mathcal{L}(k), \mathcal{C}|X)$ , defined as

$$\chi^2(k) = N \left( -1 + \sum_{i=1}^{\ell} \sum_{j=1}^k \frac{n_{ij}^2}{n_{i.} n_{.j}} \right) \quad (2)$$

is considered. It is straightforward to prove that

$$\max_{X, \mathcal{L}(k), \mathcal{C}} \chi^2(k) = N(\min\{\ell, k\} - 1) \quad (3)$$

Hence, the Ameva coefficient,  $Ameva(k) \stackrel{def}{=} Ameva(\mathcal{L}(k), \mathcal{C}|X)$ , is defined as follows:

$$Ameva(k) = \frac{\chi^2(k)}{k(\ell - 1)} \quad (4)$$

for  $k, \ell \geq 2$ . The Ameva criterion has the following properties:

- The minimum value of  $Ameva(k)$  is 0 and when this value is achieved then both discrete attributes  $\mathcal{C}$  and  $\mathcal{L}(k)$  are statistically independent and viceversa.

- The maximum value of  $Ameva(k)$  indicates the best correlation between class labels and discrete intervals. If  $k \geq \ell$  then, for all  $x \in C_i$  a unique  $j_0$  exists such that  $x \in L_{j_0}$  (remaining intervals  $(k - \ell)$  have no elements); and if  $k < \ell$  then, for all  $x \in L_j$ , a unique  $i_0$  exists such that  $x \in C_{i_0}$  (remaining classes have no elements) i.e. the highest value of the Ameva coefficient is achieved when all values within a particular interval belong to the same associated class for each interval.
- The aggregated value is divided by the number of intervals  $k$ , hence the criterion favors discretization schemes with the lowest number of intervals.
- From (3), it is followed that  $Ameva_{max}(k) \stackrel{def}{=} \max_{X, \mathcal{L}(k), \mathcal{C}} Ameva(k) = \frac{N(k-1)}{k(\ell-1)}$  if  $k < \ell$  and  $\frac{N}{k}$  otherwise. Hence,  $Ameva_{max}(k)$  is an increasing function of  $k$  if  $k \leq \ell$ , and a decreasing function of  $k$  if  $k > \ell$ . Therefore,  $\max_{k \geq 2} Ameva_{max}(k) = Ameva_{max}(\ell)$  i.e. the maximum of the Ameva coefficient is achieved in the optimal situation, it is to say, when all values of  $C_i$  are in a unique interval  $L_j$  and viceversa.

Therefore, the aim of the Ameva method is to maximize the dependence relationship between the class labels  $\mathcal{C}$  and the continuous-values attribute  $\mathcal{L}(k)$ , and at the same time to minimize the number of discrete intervals  $k$ .

## 2.2 The entropy

If  $\ell = 1$  or  $k = 1$  then it is not possible to use the Ameva method. Let us see these two cases (see Table 2 and Table 3):

Equation (2) can not be calculated using Table 2 because it

$C_i   L_j$	$L_1$	$\dots$	$L_j$	$\dots$	$L_k$	$n_{i.}$
$C_1$	$n_{11}$	$\dots$	$n_{1j}$	$\dots$	$n_{1k}$	$N$
$n_{.j}$	$n_{11}$	$\dots$	$n_{1j}$	$\dots$	$n_{1k}$	$N$

Table 2: Contingency table at first case ( $\ell = 1$ )

$C_i   L_j$	$L_1$	$n_{i.}$
$C_1$	$n_{11}$	$n_{11}$
$\vdots$	$\vdots$	$\vdots$
$C_i$	$n_{i1}$	$n_{i1}$
$\vdots$	$\vdots$	$\vdots$
$C_\ell$	$n_{\ell 1}$	$n_{\ell 1}$
$n_{.j}$	$N$	$N$

Table 3: Contingency table at second case ( $k = 1$ )

is not possible to divide by 0. Nevertheless, all the instances belong to the same class, therefore can be concluded that the dependence is maximum. In this case, let us indicate that  $A^*(1) = 1$ .

Regarding to Table 3, Ameva method can not be used because  $\chi^2(k) = 0$  and the Ameva coefficient does not give any information about the dependence. However, the dependence is not minimum and a new coefficient is necessary. By taking into account that if all instances are distributed equally in all

classes, the dependence is minimum, and if exists  $i$  such that  $n_{i1} = N$ , the dependence is maximum. Hence the following coefficient, called Entropy, is considered:

$$A(1) = 1 + \frac{1}{N \ln \ell} \sum_{i=1}^{\ell} n_{i1} \ln \left( \frac{n_{i1}}{N} \right)$$

It holds that  $0 \leq A(1) \leq 1$ , and:

- If  $A(1) = 0$ , then  $n_{i1} = \frac{N}{\ell}$  (minimum dependence).
- If  $A(1) = 1$ , then a unique  $n_{i1}$  exists that  $n_{i1} = N$  (maximum dependence).

**Note 2.1** Let us indicate these pathologic cases do not happen in a standard discretization, but it will be necessary taking into account in the presented methodology in the next section.

### 3 The methodology

Given an attribute  $X_i$  where  $i = 1, 2, \dots, s$ , the Ameva discretization algorithm is applied to this attribute so obtained intervals are considered as a new set of classes. This set of classes is denoted as follows:

$$\mathcal{C}^i = \{C_1^i, C_2^i, \dots, C_{\ell_i}^i\} \quad (5)$$

Let us consider  $X^p \subset X$  as the data subset that belongs to the class  $C_p \in \mathcal{C}^i$  where  $p = 1, 2, \dots, \ell$ . From (5), for each attribute  $X_j$  with  $j = 1, 2, \dots, s$ , a  $g_{ijp}$  value is obtained from  $\mathcal{C}^i$  as follows:

- If the  $X^p$  data subset all belong to the same class  $C^i$ , then  $g_{ijp} = A^*(1) = 1$ .
- If the subset of data belongs to different classes, then:
  - If values of the attribute  $X_j$  are always in the same interval, then  $g_{ijp} = A(1)$ .
  - If values of the attribute  $X_j$  are not always in the same interval, then  $g_{ijp} = Ameva_N(\ell_i)$ , where  $Ameva_N(\ell_i)$  is defined as follows:

$$Ameva_N(\ell_i) = \frac{\ell'_i}{N_p} Ameva(\ell_i)$$

provide that  $N_p$  is the number of instances of the class  $X^p$  and  $\ell'_i$  is the number of intervals of the attribute  $X_i$  for which there is at least one value in the data subset.

**Note 3.1** This new Ameva coefficient is chosen in order to obtain a normalized value  $0 \leq Ameva_N(\ell_i) \leq 1$  as same as  $A(1)$ .

Furthermore, it is straightforward to prove that if  $i = j$  for  $i = 1, 2, \dots, s$ , then  $g_{iip} = 1$ , for all  $p = 1, 2, \dots, \ell$ .

Given  $i, j = 1, 2, \dots, s$ , a  $g_{ij}$  value can be obtained applying this methodology for all class  $C_p \in \mathcal{C}$  ( $p = 1, 2, \dots, \ell$ ), and by considering different statistics as follows:

- The minimum  $g_{ij}^{min} = \min_p g_{ijp}$ .
- The geometric mean  $g_{ij}^{geo} = \sqrt[\ell]{\prod_{p=1}^{\ell} g_{ijp}}$ .

- The arithmetic mean  $g_{ij}^{arit} = \frac{1}{\ell} \sum_{p=1}^{\ell} g_{ijp}$ .
- The maximum  $g_{ij}^{max} = \max_p g_{ijp}$ .

It is well-known that the following relationship is holded:

$$g_{ij}^{min} \leq g_{ij}^{geo} \leq g_{ij}^{arit} \leq g_{ij}^{max}$$

The main properties of the matrix  $G = (g_{ij})$ , that is,

$$G = \begin{pmatrix} 1 & g_{12} & \cdots & g_{1s} \\ g_{21} & 1 & \cdots & g_{2s} \\ \vdots & \vdots & \ddots & \vdots \\ g_{s1} & g_{s2} & \cdots & 1 \end{pmatrix}$$

are the following: i) it is square but non symmetric matrix; ii) the values of the main diagonal are 1; and iii)  $0 \leq g_{ij}, g_{ji} \leq 1$ .

From the  $G$  matrix, a method of generating rules of dependence between attributes can be defined. For example, a possible rule is the next: given a threshold value,  $U$ , if  $\max(G_{ij}, G_{ji}) > U$  and  $i < j$ , then the  $X_j$  variable is eliminated. Let us illustrate it with an example in the next section.

### 4 A toy example

Let us consider the Iris Plants Database<sup>1</sup> from UCI Repository which is perhaps the best known database to be found in the pattern recognition literature. This data set is considered due to its simplicity since this methodology is not completely defined yet.

The data set contains four attributes (sepal length, sepal width, petal length and petal width) and three classes (Setosa, Versicolour and Virginica) of 50 instances each, where each class refers to a type of iris plant. One class is linearly separable from the other two; the latter are not linearly separable from each other.

The matrices generated by the presented methodology in this paper are:

$$G_{Iris}^{min} = \begin{pmatrix} 1 & 0.4898 & 0.6667 & 0.0998 \\ 0.3265 & 1 & 0.035 & 0.093 \\ 0.028 & 0.0586 & 1 & 0.0303 \\ 0.0545 & 0.0998 & 0.0836 & 1 \end{pmatrix} \quad (6)$$

$$G_{Iris}^{geo} = \begin{pmatrix} 1 & 0.7883 & 0.8736 & 0.4638 \\ 0.6886 & 1 & 0.3271 & 0.453 \\ 0.1727 & 0.2674 & 1 & 0.1293 \\ 0.1573 & 0.3222 & 0.2244 & 1 \end{pmatrix} \quad (7)$$

$$G_{Iris}^{arit} = \begin{pmatrix} 1 & 0.8299 & 0.8889 & 0.6999 \\ 0.7755 & 1 & 0.6783 & 0.6977 \\ 0.4039 & 0.4617 & 1 & 0.3672 \\ 0.3753 & 0.4783 & 0.4063 & 1 \end{pmatrix} \quad (8)$$

$$G_{Iris}^{max} = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{pmatrix} \quad (9)$$

This result shows that it is possible to determine the dependence of attributes of a data set from the Ameva discretization

<sup>1</sup> Available at <http://archive.ics.uci.edu/ml/datasets/Iris>

algorithm and the adjustments to resolve the inconsistencies outlined above with the entropy.

The coefficients in the minimum matrix (6) determine the lowest coefficients of dependence between two attributes. These coefficients provide information about there is a class for which the two attributes have less dependency. If these values are high, it is possible to conclude that the dependence between two attributes is high. Therefore, these coefficients are a minimum threshold for each pair of attributes.

A similar conclusion can be obtained from the maximum matrix (9). The coefficients provide information about there is a class for which the two attributes have a high dependence. In this case, these coefficients are the maximum threshold values for each pair of attributes.

Given a data set, the best result is achieved when the maximum and minimum matrix are the same. In this case, all the attributes are the same dependence with other regardless of the original class. Thus, there is only one matrix for generate the discrimination rules.

The arithmetic mean (8) and the geometric mean matrix (7) represent a global value of dependency. While the geometric mean matrix rewards the worst situations about a class, leading to a low value on the global coefficient, the arithmetic mean matrix balances the values of the coefficients.

A possible interpretation to determine which attributes are dependent of each other is to establish a threshold value. From this limit, two attributes are dependent if the average of the coefficients  $g_{ij}$  and  $g_{ji}$  of the arithmetic mean matrix is greater than or equal to this value.

In this case, the threshold value of 0.75 is established to check which attributes are dependents. The pair  $g_{ij}, g_{ji}$  that reaches this threshold is  $G_{12}, G_{21}$  because the arithmetic mean of  $G_{12}$  and  $G_{21}$  is greater than 0.75. It is necessary indicate that the sepal length and the sepal width features are the first and second attributes in the experiment.

Thus, in order to carry out a classification problem can be declared that the  $X_1$  and  $X_2$  features are similar. Let us see this affirmation by using as classification algorithm the Support Vector Machine (SVM) [1].

A performance for the 1-v-r SVM, in the form of accuracy rate, has been evaluated on models using the Gaussian kernel with  $\sigma = 1$ , and  $C = 1$ . The criteria employed to estimate the generalized accuracy is the 5-fold cross-validation on the whole set of training data. This procedure is repeated 120 times in order to ensure good statistical behavior. The obtained results are:

- With all features, the accuracy rate is 0.9320.
- Without the sepal length feature, the accuracy rate is 0.9341.
- Without the sepal width feature, the accuracy rate is 0.9667.

Furthermore to check that the accuracy rate is not less when a feature is eliminated, the methodology has discovered that these features introduce noise in the classification problem when both are used at the same time because the results are improved without the second feature.

## 5 Conclusions and future work

We have studied a method of discretization, Ameva, whose objective is to maximize the dependence between the intervals that divide the values of an attribute and the classes to which they belong, providing at the same time the minimum number of intervals.

After that, we have developed a methodology to reduce the number of features of a data set based on the dependence between them. To the best of knowledge, there are not existing researches that directly address the problem to reduce the number of features using a similar approach to ours.

This development is based on taking advantage of Ameva discretization algorithm. Thus, a new coefficient has been developed to determine the dependence between features. Hence, we have reduced the number of values of features and the number of features from a qualitative reasoning.

To test the development of the methodology, it has been applied to a well-known data set for obtain the dependent relationship between their features. Nevertheless, we think that this approach can be satisfactorily apply in this area when the data set has a lot of instances and features, and one of these features determines the class which each instance belongs. Another data sets must fulfill these characteristics.

Finally, after applying the discrimination of features obtained in the methodology, the modified data set has been carried out for the classification tests to verify the effectiveness of the methodology.

The next step to complement this development is the design of an automatic method of creation of feature discrimination rules. Subsequently, we must define some improvements in this methodology to automatically know the dependence between features without setting manually a threshold value.

## Acknowledgments

This research is supported by the Spanish Ministry of Science and Innovation R&D project ARTEMISA (TIN2009-14378-C02-01).

## References

- [1] L. González, C. Angulo, F. Velasco, and A. Catala. Dual unification of bi-class support vector machine formulations. *Pattern recognition*, 39(7):1325–1332, 2006.
- [2] Q. Wang, G.M. Garrity, J.M. Tiedje, and J.R. Cole. Naive Bayesian classifier for rapid assignment of rRNA sequences into the new bacterial taxonomy. *Applied and environmental microbiology*, 73(16):5261–5267, 2007.
- [3] M. Govindarajan. Text Mining Technique for Data Mining Application. *Proceedings of World Academy of Science, Engineering and Technology*, 26:544–549, 2007.
- [4] R. Entezari-Maleki, S.M. Iranmanesh, and B. Minaei-Bidgoli. An Experimental Investigation of the Effect of Discrete Attributes on the Precision of classification Methods. *World Applied Sciences Journal*, 7:216–223, 2009.

- [5] L. Gonzalez-Abril, FJ Cuberos, F. Velasco, and JA Ortega. Ameva: An autonomous discretization algorithm. *Expert Systems with Applications*, 36(3):5327–5332, 2009.
- [6] M. Boulle. Khiops: A statistical discretization method of continuous attributes. *Machine Learning*, 55(1):53–69, 2004.
- [7] L.A. Kurgan and K.J. Cios. CAIM discretization algorithm. *IEEE Transactions on Knowledge and Data Engineering*, 16(2):145–153, 2004.
- [8] F.J. Ruiz, C. Angulo, N. Agell, X. Rovira, M. Sánchez, and F. Prats. A Discretization Process in Accordance with a Qualitative Ordered Output. *Proceeding of the 2005 conference on Artificial Intelligence Research and Development*, 131:273–280, 2005.
- [9] R.P. Li and Z.O. Wang. An entropy-based discretization method for classification rules with inconsistency checking. 1:243–246, 2002.
- [10] L. Gonzalez-Abril, F. Velasco, JA Ortega, and FJ Cuberos. A new approach to qualitative learning in time series. *Expert Systems with Applications*, 36(6):9924–9927, 2009.
- [11] I. Guyon and A. Elisseeff. An introduction to variable and feature selection. *The Journal of Machine Learning Research*, 3:1157–1182, 2003.
- [12] G. John, R. Kohavi, and K. Pflieger. Irrelevant features and the subset selection problem. *Machine Learning Conference Proceedings*, pages 121–129, 1994.
- [13] J. Yang and V. Honavar. Feature subset selection using a genetic algorithm. *Intelligent Systems and Their Applications, IEEE*, 13(2):44–49, 1998.
- [14] KM Faraoun and A. Rabhi. Data dimensionality reduction based on genetic selection of feature subsets. *INFOCOM - Journal of Computer Science*, 6(2):9–19, 2007.
- [15] L. Rocchi, L. Chiari, and A. Cappello. Feature selection of stabilometric parameters based on principal component analysis. *Medical and Biological Engineering and Computing*, 42(1):71–79, 2004.
- [16] Nitin Khosla. *Dimensionality Reduction Using Factor Analysis*. PhD thesis, 2006.