

A TUNING AID FOR DISCRETIZATION IN RULE INDUCTION

M. POSTEMA, X. WU, T. MENZIES

*Department of Software Development, Monash University, 900 Dandenong Rd,
Melbourne, Vic 3145, AUSTRALIA*

This paper examines where a tuning aid can be useful to help discretization of numerical attributes in rule induction, and subsequently improve deduction of induction results. Different discretization methods use different strategies to set up the borders for continuous attributes. They mostly incorporate class supervision to define the discretization borders. The tuning aid we present uses an unsupervised method to define the intervals at induction time. We then supervise the learning process, by comparing the performance in terms of predictive accuracy with the information gain based discretization methods implemented in HCV (Version 2.0) and C4.5.

1 Introduction

In the context of rule induction and decision tree construction, dealing with a continuous attribute means discretization of the numerical attribute into a number of intervals. The discretized intervals can be treated in a similar way to nominal values during induction and deduction. The essential aspect of discretization is to find the right places to set up interval borders. In supervised discretization methods, such as the information gain based methods implemented in C4.5⁹ and HCV (Version 2.0)¹¹, the class information of examples in a training set is used. In unsupervised (or class-blind) discretization methods, such as equal-width discretization and equal-frequency discretization², training examples are grouped into intervals without taking into account the respective classes of the training examples. Discretization can be performed at induction time (such as in C4.5) or before induction takes place (see³ and⁷).

Among various discretization methods, the information gain based methods have been widely used and cited^{3,12}. C4.5 provides only binary discretization at induction time based on an information gain approach for real-valued attributes. In HCV (Version 2.0), an information gain based method is the default discretization method for processing numerical attributes before induction takes place. In this paper, we present an unsupervised tuning aid, based on domain distribution, for discretization of real-valued attributes. The method has been incorporated into HCV (Version 2.0) as a counterpart to the information gain method, and we compare its performance with the information gain method also implemented in HCV (Version 2.0). The results of C4.5 with its binary discretization method are also included for comparison,

although this inclusion might not be relevant because C4.5 and HCV (Version 2.0) have different induction strategies.

The rest of this paper is organized as follows. Section 2 outlines HCV (Version 2.0) and the information gain heuristic, Section 3 discusses the fine tuning design, Section 4 details the experiments, and we follow this with results and finally a conclusion.

2 HCV (Version 2.0) and the Information Gain Heuristic

2.1 HCV (Version 2.0)

The HCV algorithm¹⁰ is a representative of the extension matrix based family of attribute-based induction algorithms, originating with J.R. Hong's AE1⁶. By dividing the positive examples (PE) of a specific concept in a given example set into intersecting groups and adopting a set of strategies to find a heuristic conjunctive rule in each group which covers all the group's positive examples and none of the negative examples (NE), HCV can find a rule in the form of variable-valued logic for the concept in low-order polynomial time. If there exists at least one conjunctive rule in a given training example set for PE against NE, the rule produced by HCV must be a conjunctive one. The rules in variable-valued logic generated by HCV have been shown¹¹ empirically to be more compact than the decision trees or their equivalent decision rules produced by the ID3 algorithm (the best-known induction algorithm to date) and its successors (e.g., C4.5) in terms of the numbers of conjunctive rules and conjunctions.

The HCV (Version 2.0) software is a C++ implementation of the HCV algorithm. In this implementation, HCV can work with noisy and real-valued domains as well as nominal and noise-free databases. It also provides a set of deduction facilities for the user to test the accuracy of the produced rules on test examples. The detailed description of the software is included in¹¹.

In addition to a set of discretization facilities, such as the information gain heuristic outlined in the next section, HCV (Version 2.0) permits the user to specify their own discretization of real-valued attributes by providing a set of intervals in the structure file, which specifies the attributes (with their order and value domains) and classes used in the data files. This is a very useful way for integrating domain specific information. The tuning aid designed in Section 3 will start with this facility.

2.2 The Information Gain Heuristic for Discretization

When the examples in a training set have taken values of x_1, \dots, x_n in ascending order on a continuous attribute, we can use the information gain heuristic adopted in ID3⁸ to find a most informative border to split the value domain of the continuous attribute. Fayyad and Irani⁴ have shown that the maximum information gain by the heuristic is always achieved at a cut point (say, the mid-point) between the values taken by two examples of different classes.

The information gain heuristic is adopted in HCV (Version 2.0) in the following way. Each $x = (x_i + x_{i+1})/2$ ($i = 1, \dots, n - 1$) is a possible cut point if x_i and x_{i+1} have been taken by examples of different classes in the training set. Use the information gain heuristic to check each of the possible cut points and find the best split point. Run the same process on the left and right halves of the splitting to split them further. The number of intervals produced this way may be very large if the attribute is not very informative. Catlett¹ has proposed some criteria to stop the recursive splitting which have been adopted in HCV (Version 2.0):

- Stop if the information gain on all cut points is the same,
- Stop if the number of examples to split is less than a certain number (*e.g.* fourteen in HCV (Version 2.0)), and
- Limit the number of intervals to be produced to a certain number (*e.g.* eight in HCV (Version 2.0)).

In C4.5⁹, the information gain approach is revised in the following ways. Firstly, each of the possible cut points is not the midpoint between the two nearest values, but rather the greatest value in the entire training set that does not exceed the midpoint. This ensures that all border values occur in the training data. Each border value in this case is not necessarily the same as the lower of the two neighbouring values since all training examples are examined for the selection. Secondly, C4.5 adopts the information gain ratio rather than the information gain heuristic. Finally, C4.5 does binarization of continuous attributes, which means only one interval border is found for each continuous attribute at each decision node.

3 A Fine Tuning Aid to HCV (Version 2.0)

3.1 Related Work

The tuning aid that we present in Section 3.2 is similar to *Adaptive Quantizers* discussed in Dougherty³. The Adaptive Quantizers method begins with

a binary equal width partitioning of the continuous feature. Induction is performed and tested for predictive accuracy. The interval with the lowest predictive accuracy is split into two equal width partitions. The induction and evaluation processes are then repeated until some performance criteria is met. This method appears to overcome some of the limitations of unsupervised discretization, but comes with a high computational cost. This is due to numerous runs of the rule induction process. The method also assumes that a high level of accuracy can be achieved.

The minimal entropy heuristic^{1,3} uses class information entropy of candidate partitions to select the intervals.

The method we present uses some of the above mentioned algorithms slightly differently. Based on the HCV (Version 2.0) facility for user-specified intervals (see Section 2.1), the tuning aid method for discretization is designed in Section 3.2.

3.2 Fine Tuning Design

We use an unsupervised learning method that firstly partitions the continuous attribute into ten equal width intervals. We then find three intervals based on the data, not the class distribution. These intervals are then incorporated into the induction algorithm. If the performance decreases (compared to the default) we don't do any further processing. Whilst the computational cost increases (refer Section 4.2), this is marginal when compared to the improvement of results. The hypothesis and tuning aid are presented in Sections 3.2.1 and 3.2.2 respectively.

Hypothesis

Initial experiments were devised based on the attribute information available in the dictionary files. As is standard with data sets from the University of California at Irvine Repository of Machine Learning Databases and Domain Theories, a dictionary file specifies the attributes (with their order and value domains) and classes used in the training and test data. From here a hypothesis was formed that the overall data distribution could be analyzed to determine intervals in the *low*, *average* and *high* range. This was facilitated by plotting a histogram of the data in the training files into ten equal intervals. These intervals were then analyzed to find the following three intervals:

- *Low* which contained approximately 25% of the data
- *Average* which contained approximately 50% of the data

- *High* which contained the remaining data

Since it is not possible to use this method when a large percentage (say 70%) of the data lies at either extremes of the range of values, a second method is introduced which includes overlapping class intervals. For example, if class 1 has values ranging 0-100 and class 2 has values ranging 85-150, the intervals used are 0-100 and 85-150. This method had limited use and success.

Tuning Aid

1. Run the HCV (Version 2.0) software on a database to obtain baseline results. The database should consist of at least one continuous attribute.
2. Select a continuous attribute in the database.
3. Check the overall data distribution on the selected attribute.
4. Divide the continuous attribute into the following three intervals
 - low (approximately 25% of the data)
 - average (approximately 50% of the data)
 - high (approximately 25% of the data)

If the above fails attempt to determine overlapping class intervals (as mentioned in Section 3.2.1).

5. If intervals were found in step 4, run the HCV (Version 2.0) software with these intervals and obtain results, else go to step 7.
6. If the results are improved or no worse than the baseline ones, keep them, else discard these intervals.
7. If no more continuous attributes, STOP.
8. Select another continuous attribute.
9. Go to step 3.

The second method discussed in Section 3.2.1 and Step 4 above deals with domains with some degree of fuzziness. For example, we could classify someone who is 170cm in height as ‘Normal Height’ or ‘Tall Height’ depending on our outlook. Instead of discretizing the domain into fixed intervals, some sort of curve could be applied to give a fuzzy border. Thus a value could belong to more than one interval at the same time. Results can be classified in terms

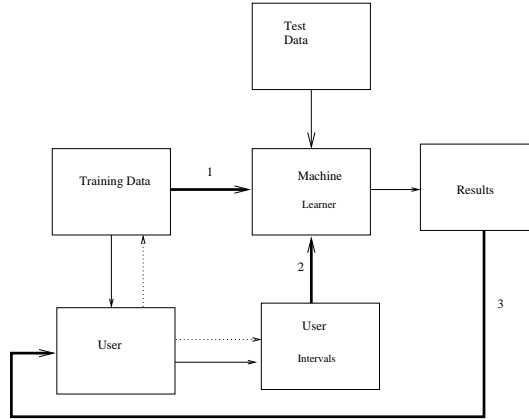


Figure 1: The Machine Learning Cycle

of degrees of membership in each interval. Discretizing a continuous domain in this way is termed *fuzzy matching*¹³. Alternatively, the number of intervals and their values could also be modified, which is the principle used in fuzzy expert systems. This above method combines probabilistic analysis (refer to Step 4) and fuzzy logic.

The method based on the information in the dictionary file and the data in the training files can be related to the probability surveys between domain experts when designing a fuzzy expert system⁵. A summary of the tuning aid method can be included into the machine learner as shown by Figure 1, where the user can specify the intervals from the training examples. Following the heavy arrowed lines, we can see:

1. Training and test data are fed into the machine learner as a baseline.
2. The user specifies intervals, feeds these into the machine learner and runs the process again.
3. The deduction of induction results on test data are then compared to determine which method gave the best results.

The dotted lines indicate the user decision.

Table 1: Success Analysis of Results

<i>Domain</i>	<i>Success Analysis</i>
Wine	$\oplus \ominus \oplus \oplus \oplus \oplus \ominus \ominus \ominus$
Bupa	$\oplus \oplus \ominus \ominus \ominus \ominus$
Labor Neg	$\oplus \oplus \ominus \ominus \ominus \ominus$
Swiss 5	$\oplus \ominus \ominus \ominus$
Cleveland 2	\oplus
Cleveland 5	\odot
Va 5	\oplus
Va 2	\ominus
Crx	\oplus
Imports 85	$\odot \odot \ominus \odot \ominus \ominus \ominus \oplus \odot \ominus \ominus \oplus \ominus \ominus$
Pima	$\ominus \ominus \ominus \ominus \ominus \oplus$
Glass (without ID)	$\ominus \ominus \ominus \ominus \ominus \ominus$

Legend

- \odot indicates same accuracy (intervals were included in subsequent experiments).
- \ominus indicates decreased accuracy (subsequently discarded).
- \oplus indicates improved result.

4 Experiments

The experiments were performed with data sets from the University of California at Irvine Repository of Machine Learning Databases and Domain Theories. The inputs for the machine learners include the *dictionary*, *training* and *testing* files, which describe the attributes, and include training and testing examples respectively.

4.1 Baseline

The experimental results can be compared to those obtained by Wu et al.¹⁴. The best results in terms of predictive accuracy for each problem is given in *bold* in Table 2.

4.2 Results

Table 1 analyses the success of using the data distribution method defined in Section 3.2 to specify the intervals in continuous domains. One interesting observation is that the probabilistic method is actually a search for attributes that can be discretized usefully. In the twelve data sets, 16 of the 55 continuous

Table 2: Summary and Comparison of Results

Domain	Accuracy (%)		
	C4.5	HCV (Ver2.0)	HCV (Tuned)
Bupa	73.7	57.6	65.3
Cleveland 5	47.3	56.0	56.0
Cleveland 2	68.4	78.0	82.4
Crx	79.5	82.5	85.0
Imports 85	64.4	62.7	66.1
Labor Neg	82.4	76.5	82.4
Pima	75.1	73.9	75.1
Swiss 5	31.2	28.1	34.4
Va 5	26.8	26.8	32.4
Va 2	66.7	78.9	70.4
Wine	90.4	90.4	98.1
Glass (without ID)	64.6	72.3	69.2

attributes experimented with give improved results, 34 give worse results and 5 make no difference.

The induction algorithm is run $N+1$ times where N represents the number of continuous attributes, thereby increasing the computational cost by N times.

The experiments including the probabilistic method are recorded as HCV (Tuned) in Table 2. The overall success in terms of accuracy is high if we consider that results from HCV (Tuned) have the highest accuracy on 7 out of the 12 data sets tested, and matched C4.5 in another two.

5 Conclusion

The success of fuzzy expert systems in the area of domain specific controllers has been one of the stimulations for further research work in improving machine learning results. This paper has presented a tuning aid method, which combines probability analysis and fuzzy logic, for discretization of continuous attributes, and has compared its results with the information gain based discretization methods implemented in HCV (Version 2.0) and C4.5 on various problem domains. Experimental results have shown that the tuning aid method has given the highest accuracy on 7 out of the 12 data sets tested, with equal results in two other data sets. Based on the analysis of fuzzy expert systems, future work will include optimizing the number of intervals and incorporating the tuning aid into other machine learners, such as C4.5.

References

1. J. Catlett. On Changing Continuous Attributes into Ordered Discrete Attributes. *Proceedings of the 1991 European Working Session on Learning*, 1991.
2. D. Chiu, A. Wong and B. Cheung. Information Discovery through Hierarchical Maximum Entropy Discretization and Synthesis. In Piatetsky-Shapiro, G., and W.J. Frowley, editors, *Knowledge Discovery in Databases*. MIT Press, 1991.
3. J. Dougherty, R. Kohavi and M. Sahami. Supervised and Unsupervised Discretization of Continuous Features. *Proceedings of the 12th International Conference on Machine Learning*, 194–202, 1995.
4. U.M. Fayyad and K.B. Irani. On the Handling of Continuous-Valued Attributes in Decision Tree Generation. *Machine Learning*, **8**(1992), 87–102.
5. L.O. Hall and A. Kandel. The Evolution of Expert Systems. In F. Aminzadeh and M. Jamshidi, editors, *Soft Computing: Fuzzy Logic, Neural Networks, and Distributed Artificial Intelligence*. Prentice Hall, New Jersey USA, 1994.
6. J.R. Hong. AE1: An Extension Matrix Approximate Method for the General Covering Problem. *International Journal of Computer and Information Sciences*, **4**(1985): 421–437.
7. B. Pfahringer. Compression-Based Discretization of Continuous Attributes, *Proceedings of the 12th International Conference on Machine Learning*, 456–463, 1995.
8. J.R. Quinlan. Induction of Decision Trees. *Machine Learning*, **1**(1986), 81–106.
9. J.R. Quinlan. *C4.5: Programs for Machine Learning*. Morgan Kaufmann Publishers, 1993.
10. X. Wu. The HCV Induction Algorithm. In S.C. Kwasny and J.F. Buck, editors, *Proceedings of the 21st ACM Computer Science Conference*, 168–175. ACM Press, USA, 1993.
11. X. Wu. *Knowledge Acquisition from Databases*. Ablex, USA, 1995.
12. X. Wu. A Bayesian Discretizer for Real-Valued Attributes. *The Computer Journal*, **39**(1996).
13. X. Wu and P. Máhlén. Fuzzy Interpretation of Induction Results. *Proc. of the 1995 International Conference on Knowledge Discovery and Data Mining (KDD-95)*, Montreal, Canada, August 20–21, 1995, 325–330.
14. X. Wu, J. Krisár, and P. Máhlén. Noise Handling with Extension Matrices. *International Journal of Artificial Intelligence Tools*, **5**(1996), 81–97.