(I)

95　10　11

(I)

( )

95 10 8

into account the effect of degrees of freedom. However, the classification with a controlled degree of uncertainty, or a misclassification error, is outside the realm of RST. This algorithm also ignores the effect of variance in the two merged intervals. In this study we propose a new algorithm, named the extended Chi2 algorithm, to overcome these two drawbacks. By running the software of See5, our proposed algorithm possesses a better performance than the original and modified Chi2 algorithms.

Chi2

(extended Chi2 algorithm)

(Chi2 , Chi2 )

**Abstract**

The Variable Precision Rough Sets (VPRS) model is a powerful tool for data mining, as it has been widely applied to acquire knowledge. Despite its diverse applications in many domains, the VPRS model unfortunately cannot be applied to real world classification tasks involving continuous attributes. This requires a discretization method to pre-process the data. Discretization is an effective technique to deal with continuous attributes for data mining, especially for the classification problem. The modified Chi2 algorithm is one of the modifications to the Chi2 algorithm, replacing the inconsistency check in the Chi2 algorithm by using the quality of approximation, coined from the Rough Sets Theory (RST), in which it takes

A number of methods based on the entropy measure establish a strong group of works in the discretization domain. This concept uses class entropy as a criterion to evaluate a list of best cuts, which together with the attribute domain induce the desired intervals (Nguyen, 1998).

The ChiMerge algorithm introduced by Kerber (1992) is a supervised global discretization method. The user has to provide several parameters such as the significance level $\alpha$, and the maximal intervals and minimal intervals during the application of this algorithm. ChiMerge requires $\alpha$ to be specified. Nevertheless, too big or too small a $\alpha$ will over-discretize or under-discretize an attribute. Liu, et al. (1997) proposed a Chi2 algorithm that uses a ChiMerge algorithm as a basis, whereby the Chi2 algorithm improves the ChiMerge algorithm in that the value of $\alpha$ is calculated based on the training data itself.

Tay, et al. (2002) indicated that although the Chi2 algorithm automates the ChiMerge

algorithm by calculating a significance value $\alpha$ based on the training data set, it still has two drawbacks: (1) the Chi2 algorithm requires the user to provide an inconsistency rate to stop the merging procedure. This is unreasonable since an inappropriate threshold will result in over-merging. (2) This merging criterion does not consider the degrees of freedom, but rather only the fixed degrees of freedom (the classes' number minus one). According to the statistical point of view, this is inaccurate (Montgomery, et al. 1999), since the power of a statistical test is affected by the degrees of freedom of a test. They utilize the quality of approximation to replace the inconsistency checking of the Chi2 algorithm and consider the degrees of freedom of each two adjacent intervals, in which the two adjacent intervals when it has a maximal difference in the calculated $\chi^2$ value and the threshold should be merged first.

The modified Chi2 algorithm introduced by Shen et al. (2001) can be sectioned into two phases: The first phase of the modified Chi2 algorithm can be regarded as a generalization version of the ChiMerge algorithm. Instead of specifying a $\chi^2$ threshold, the modified Chi2 algorithm provides a wrapping that automatically increments the $\chi^2$ threshold (decreasing the significant level $\alpha$). A consistency check is used as a stopping criterion to make sure that the modified Chi2 algorithm automatically determines a proper $\chi^2$ threshold while still keeping the fidelity of the original data.

The second phase is a finer process of the first phase, beginning with the significant level $\alpha_0$ determined in the first phase, where each attribute $i$ is associated with a sigLvl[$i$] and they take turns for merging. A consistency check is conducted after each attribute's merging. If the inconsistency rate does not exceed the pre-defined inconsistency rate ($\delta$), then sigLvl[$i$] is decreased for attribute $i$'s next round of merging. Otherwise, the attribute $i$ will not be involved in further merging. This process is repeated until no attribute's value can be merged.

In the modified Chi2 algorithm, inconsistency checking (InConCheck (data) $<\delta$) of the original Chi2 algorithm is replaced by the quality of approximation $L_c$ after each step of discretization ($L_{c-discretized} \leq L_{c-original}$). This inconsistency rate is utilized as the termination criterion. The quality of approximation coined from the Rough Sets Theory is defined as follows:

$$L_c = \frac{\sum card(\underline{B}X_i)}{card(U)}, \qquad (2.1)$$

where $U$ is the set of all objects of the data set:

$X$ can be any subset of $U$;

$\underline{B}X$ is the lower approximation of $X$ in $B$ ($B \subseteq A$);

$A$ is the set of attributes.

The card denotes set cardinality.

The merge criterion of the original Chi2 algorithm does not consider the degrees of freedom, as it only used the fixed degrees of freedom (the classes' number minus one). The original Chi2 algorithm merges the pair of adjacent intervals with the lowest $x^2$ value being the critical value. The merge criterion of modified Chi2 considers the degrees of freedom of each of the two adjacent intervals. When two adjacent intervals have a maximal difference in the calculated $\chi^2$ value, the threshold should be merged first.

**The extended Chi2 algorithm**
Step 1. Initialize:
Set the significant level as $\alpha = 0.5$;
calculate the pre-defined

inconsistency rate $\xi$ .

Step 2. Calculate the chi-square value:

For each numeric attribute, sort data on the attribute and use formula (3.2) to compute the $x^2$ value.

Step 3. Merge:

For a comparison, compute the $x^2$ value and corresponding threshold; merge the adjacent two intervals which have the maximal normalized difference and the computed $x^2$ value is smaller than the corresponding threshold. If no two adjacent intervals satisfy this condition, then go to Step 5.

Step 4. Check inconsistency rate for merger:

Check the merged inconsistency rate, and if the merged inconsistency rate exceeds the pre-defined inconsistency rate, then discard the merger. Go to step 5. Otherwise, go to step 2.

Step 5. Decrease the significance level:

Decrease $\alpha \to \alpha_0$ .

Step 6. Calculate finer the chi-square value:

For each numeric attribute, sort data on the attribute and use formula (3.2) to compute the $x^2$ value.

Step 7. Finer merge:

For a comparison, compute the $x^2$ value and corresponding threshold; merge the adjacent two intervals which have a maximal normalize difference and the computed $x^2$ value is smaller than the corresponding threshold. If no two adjacent intervals satisfy this condition, then go to Step 9.

Step 8. Check the inconsistency rate much finer for a merger:

Check the merged inconsistency rate; if the merged inconsistency rate exceeds the pre-defined inconsistency rate, then discard the merger. Go to step 9. Otherwise, go to step 6.

Step 9. Decrease finer the significance level:

Decrease the significance level; then stop.

The formula for computing the $\chi^2$ value

is: $\chi^2 = \sum_{i=1}^{n}\sum_{j=1}^{k}\frac{(A_{ij}-E_{ij})^2}{E_{ij}}$ , (3.2)

where $n = 2$ ;

$k =$ number of classes;

$A_{ij} =$ number of objects in the $i$th interval, $j$th class;

$R_i =$ number of objects in the $i$th interval $= \sum_{j=1}^{k} A_{ij}$ ;

$C_j =$ number of objects in the $j$th class $= \sum_{i=1}^{n} A_{ij}$ ;

$N =$ total number of objects $= \sum_{i=1}^{n} R_i$ ;

$E_{ij} =$ expected frequency of

$A_{ij} = \dfrac{R_i * C_j}{N}$ .

If either $R_i$ or $C_j$ is 0, then $E_{ij}$ is set to 0.1. The degrees of freedom of the $\chi^2$ statistic are one less than the number of classes

Five data sets are demonstrated to present the effectiveness of the proposed extended Chi2 algorithm. The five data sets are taken from the University of California, Irvine's repository of machine learning databases.

We ran See5 on both the original data sets and the discretized data sets. The parameters

of See5 utilize its default setting. The ten-fold cross-validation test method is applied to all data sets. The data set is divided into 10 parts of which nine parts are used as training sets and the remaining one part as the testing set. The experiments were repeated 10 times. The final predictive accuracy is taken as the average of the 10 predictive accuracy values.

T    The extended Chi2 algorithm is compared with the original Chi2 and modified Chi2 and Boolean reasoning algorithm with the predefined inconsistency rate ($\delta$) value equal to 0 in the experiment.

The discretized data sets are sent into See5. The predictive accuracy and its standard deviation of these methods are listed in Table 1. From Table 1, we know that the predictive accuracy of the extended Chi2 algorithm outperforms other discretization algorithms.

Table1.    The Predictive Accuracy Using See5 With the Discretization Algorithm

| Data Set | See5 | | | | |
|---|---|---|---|---|---|
| | Continuous | Original Chi2 Algorithm | Modified Chi2 Algorithm | Extended Chi2 Algorithm | Boolean Reasoning Algorithm |
| Bupa | $67.5 \pm 2.4\%$ | $65.2 \pm 3.2\%$ | $67.5 \pm 1.9\%$ | $68.4 \pm 2.7\%$ | $68.1 \pm 2.3\%$ |
| Glass | $68.6 \pm 2.5\%$ | $93.1 \pm 2.1\%$ | $93.4 \pm 2.3\%$ | $93.5 \pm 1.3\%$ | $71.9 \pm 2.8\%$ |
| Iris | $94.0 \pm 2.1\%$ | $94.0 \pm 2.1\%$ | $93.3 \pm 2.2\%$ | $94.0 \pm 2.1\%$ | $96.0 \pm 1.8\%$ |
| Breast Cancer | $94.9 \pm 0.8\%$ | $95.5 \pm 1.0\%$ | $96.0 \pm 0.9\%$ | $96.5 \pm 0.8\%$ | $95.2 \pm 0.8\%$ |
| Heart dissease | $51.9 \pm 1.4\%$ | $52.5 \pm 2.3\%$ | $53.2 \pm 2.7\%$ | $54.2 \pm 1.7\%$ | $55.9 \pm 2.6\%$ |

Many classification algorithms developed in the data mining community can only acquire knowledge on the nominal attributes' data sets. However, many real world classification tasks exist that involve continuous attributes, such that these algorithms cannot be applied unless the continuous attributes are discretized. The VPRS model is a powerful mathematical tool for data analysis and knowledge discovery from inconsistent and ambiguous data. It cannot be applied to extract rules from the continuous attributes unless they are first discretized.

In this study we propose an extended Chi2 algorithm that determines the pre-defined misclassification rate ($\delta$) from the data itself. We also consider the effect of variance in the two adjacent intervals. With these modifications, the extended Chi2 algorithm not only handles misclassified or uncertain data, but also becomes a completely automated discretization method and its predictive accuracy is better than the original Chi2 algorithm.

Five real world data set experiments were conducted to demonstrate the feasibility of the proposed algorithm. The experimental results show that our proposed algorithm could acquire a higher predicted accuracy than the original and modified Chi2 algorithm. Furthermore, the tree size is significantly smaller than using the original data with See5.

For $m$ attributes, the computational complexity of original Chi2 algorithm at phase 1 has $O(Kmn \log n)$, where $n$ is the number of objects in the dataset, and $K$ is the number of incremental steps. A similar complexity can be obtained for phase 2. Although our proposed algorithm adds one step (i.e., to select the merging intervals), it

does not increase the computational complexity as compared to the original Chi2 algorithm. The computational complexities of the original Chi2 algorithm, modified Chi2 algorithm, and our proposed algorithm are the same.

The above research results have been accepted for publication in *IEEE Transactions on Knowledge and Data Engineering.*

1. An, A., Shan, N., Chan, C., Cercone, N. and Ziarko, W., 1996, "Discovering Rules for Water Demand Prediction: An Enhanced Rough-set Approach," *Engineering Applications in Artificial Intelligence*, Vol. 9, No.6, pp. 645-653

2. Beynon, M., 2001, "Reducts within the variable precision rough sets model: A further investigation," *European Journal of Operational Research*, Vol. 134, pp. 592-605.

3. Beynon, M., 2002, "The Identification of Low-Paying Workplaces: An Analysis Using the Variable Precision Rough Sets Model," *The Third International Conference on Rough Sets and Current Trend in Computing, Lecture Notes in Artificial Intelligence Series*, Springer-Verlag, pp. 530-537.

4. Chmielewski, R. and Grzymala-Busse, W., 1996, "Global Descretization of Continuous Attributes as Preprocessing for Machine Learning," *International Journal of Approximate Reasoning*, Vol. 15, No. 4, pp. 319-331.

5. Dougherty, J., Kohavi, R. and Sahami, M., 1995, "Supervised and Unsupervised Discretization of Continuous Features," *Machine Learning: Proceedings of the Twelfth International Conference*, San Francisco, pp. 194-202.

6. Kattan, M. W. and Cooper, R. B., 1998, "The Predictive Accuracy of Computer-based Classification Decision Techniques. A Review and Research Directions," *Omega-International Journal of Management Science*, Vol. 26, No. 4, pp. 467-482.

7. Kerber, R., 1992, "ChiMerge: Discretization of Numeric Attributes," Proceeding tenth International Artificial Intelligence, pp. 123-128.

8. Li, R. P. and Wang, Z. O., 2002, "An Entropy-Based Discretization Method for Classification Rules with Inconsistency Checking," *Proceedings of the First Conference on Machine Learning and Cybemetics*, Beijing, pp. 243-246.

9. Liu, H. and Setiono, R., 1997, "Feature Selection via Discretization," *IEEE Transactions on Knowledge and Data Engineering*, Vol. 9, No. 4, pp. 642-645.

10. Montgomery, D. C. and Runger, G. C., 1999, *Applied Statistics and Probability for Engineers*, Jone Wiley & Sons.

11. Nguyen, H. S. and Nguyen, S. H., 1998, "Discretization Methods in Data Mining," *Rough Sets in Knowledge Discovery*, Physica-Verlag, Heidelberg, pp. 451-482.

12. Shen, L. and Tay, E. H., 2001, "A discretization method for Rough Sets Theory," Intelligent Data Analysis, Vol. 5, pp. 431-438.

13. Tay, E. H. and Shen, L., 2002, "A Modified Chi2 Algorithm for Discretization," *IEEE Transactions on Knowledge and Data Engineering*, Vol. 14, No. 3, pp. 666-670.

14. Ziarko, W., 1993, "Variable Precision Rough Set Model," *Journal of Computer and System Science*, Vol.46, pp. 39-59.