# *Maejo International Journal of Science and Technology*

*Full Paper*

# Enhancement of transparency and accuracy of credit scoring models through genetic fuzzy classifier

**Adel Lahsasna, Raja N. Ainon\* and  Teh Y. Wah**

Faculty of Computer Science and Information Technology, University of Malaya, Kuala Lumpur 50603, Malaysia

\*Corresponding author, e-mail: ainon@um.edu.my

**Abstract:** Credit risk evaluation systems play an important role in the financial decision-making by enabling faster credit decisions, reducing the cost of credit analysis and diminishing possible risks. Credit scoring is the most commonly used technique for evaluating the creditworthiness of the credit applicants. The credit models built with this technique should satisfy two important criteria, namely accuracy, which measures the capability of predicting the behaviour of the customers, and transparency, which reflects the ability of the model to describe the input-output relation in an understandable way. In our paper, two credit scoring models are proposed using two types of fuzzy systems, namely Takagi-Sugeno (TS) and Mamdani types**.** The accuracy and transparency of these two models have been optimised. The TS fuzzy credit scoring model is generated using subtractive clustering method while the Mamdani fuzzy system is extracted using fuzzy C-means clustering algorithm. The accuracy and transparency of the two resulting fuzzy credit scoring models are optimised using two multi-objective evolutionary techniques. The potential of the proposed modelling approaches for enhancing the transparency of the credit scoring models while maintaining the classification accuracy is illustrated using two benchmark real world data sets. The TS fuzzy system is found to be highly accurate and computationally efficient while the Mamdani fuzzy system is highly transparent, intuitive and humanly understandable.

**Keywords:**  credit scoring, fuzzy classifier, genetic algorithms, transparency

## Introduction

The global financial crisis of 2008 reveals the importance of the credit risk evaluation decisions not only on the financial institutions and banks but also on both global and local economy. Many major banks collapsed and others suffered heavy losses as a result of mortgage payment default. Hence, decision support tools that aim to enhance the manager's decision may play a valuable role in decision-making by allowing faster and more accurate decisions.

Credit scoring is the most commonly used method for evaluating the creditworthiness of the applicants. Before this method came into use, judgmental method was the only way to differentiate between the good applicants who are likely to repay their debts and the bad ones who are denied because of the high potential of defaulting on their debts. This approach to credit assessment has been criticised for being inconsistent, costly and time consuming. In recent years, credit scoring which replaced the judgmental method, aims at classifying credit applicants into bad and good customers with respect to their features such as age, income, and marital status [1].

Accuracy and transparency are two important criteria that should be satisfied by any credit scoring model. A highly accurate credit model enables correct assessment, thus avoiding any heavy losses associated with wrong predictions while transparent credit model enables financial analysts to understand the decision process.

The literature on credit scoring shows that statistical methods such as linear discriminant analysis and logistic regression are the most commonly used methods in building credit scoring models [2]. However, artificial intelligence techniques such as neural networks and genetic algorithms provide a new alternative to statistical methods in optimising non-linear, complex and real world systems [1, 3-5].

The main reason artificial intelligence techniques are seldom used in credit risk evaluation industry is the lack of explanatory capabilities of these methods. Hence, the enhancement of the transparency of the artificial-intelligence-based credit scoring model is one of the key factors of their successful deployment [6]. The main advantage of the fuzzy system is its transparency. Through the hybridisation of the transparency of the fuzzy system with the excellent learning capacity of the artificial intelligence techniques, some limitations of single-methods transparency may be overcome. Using this approach, some credit scoring models have been proposed using neuro-fuzzy [7-8] and genetic fuzzy [9] techniques to solve the transparency problem. Hoffman et al. [9] proposed a genetic fuzzy system for credit scoring and compared it with Nefclass, a neuro-fuzzy algorithm. The results showed that the performance of the genetic fuzzy algorithm is better than Nefclass [10] while the latter is more transparent. In addition, the above stated study reveals the classical trade-off between accuracy and transparency in the fuzzy systems. Hence, such problem has to be carefully addressed and balanced based on the needs and the objective of the credit scoring user.

In a recent study [3], the main soft computing methods applied in credit scoring models were surveyed. However, the multi-objective genetic algorithm, which is an efficient technique to get a maximum trade-off between conflicting objectives, has not been investigated for its handling of the accuracy and transparency trade-off in the fuzzy-based credit scoring models. The multi-objective genetic algorithm has, however, been successfully applied in the design phase of the

fuzzy-rule-based system modelling [11]. Specifically, it has been used in this phase to find an appropriate balance between transparency and complexity of the fuzzy-rule-based system [12]. Moreover, the multi-objective Pareto optimal solutions, the adopted approach in this paper, give more realistic solutions to the problem by allowing the decision-maker to choose between different solutions based on his needs and conditions. This paper aims at investigating the significance of this approach for addressing the above stated problem using two real-world data sets.

In this paper, two credit scoring models are built using Takagi-Sugeno (TS) and Mamdani fuzzy systems. Particularly, the accuracy and transparency of the resulting credit fuzzy models are enhanced using two different multi-objective genetic algorithms. To illustrate the potentiality of the proposed methods, two benchmark data sets, namely German [13] and Australian [14] credit data sets, are used. An overview of the multi-objective genetic optimisation methods is first detailed, followed by the description of the adopted methodology and the data sets used in this study.

**Multi-Objective Genetic Algorithms**

Many real-world problems have multiple conflicting objectives that should be simultaneously considered, as the optimisation of a particular solution with respect to one objective can give unacceptable results with respect to other objectives. A reasonable approach to multi-objective optimisation problem is to find a set of solutions, each of which achieves the objectives in a balanced way without being dominated by any other solution. Genetic algorithms, the meta-heuristic techniques inspired by the evolutionary biology, are well suited to this class of problems [15].

There are two approaches in multi-objective genetic algorithms optimisation. The first is to combine the various objective functions into a single function in a linear fashion using weight factors. The drawback of this approach lies in the determination of the optimal weight values that characterise the user preferences. The second approach finds the non-dominated Pareto optimal set of solutions for all optimal compromises between the conflicting objectives. It is a practical approach as the decision-maker can find solutions with different trade-off levels. A number of algorithms have been proposed [16-17] and the elitist non-dominated sorting genetic algorithm II (NSGA-II) [18] is among the well-known and most commonly used multi-objective genetic algorithms in the literature.

The NSGA-II was introduced to overcome the following drawbacks of NSGA [19]: (i) computation complexity, (ii) non-elitism approach and (iii) the need for specifying a sharing parameter. This algorithm has two features which makes it an efficient algorithm. The first one is that the fitness function of the solution is based on non-dominated ranking and a crowding measure, and the second is the elitist-generation update procedure. A non-dominated rank is assigned to each individual using the relative fitness. The concept of non-dominated solution can be defined as follows: Individual or solution 'A' dominates 'B' if the two following conditions hold:
(i)     'A' is strictly better than 'B' in at least one objective and
(ii)     'A' is no worse than 'B' in all objectives.

An outline of the elitism-preserving mechanism of NSGA-II is written as follows:

Step 1: Generate an initial population with *N* chromosomes.

Step 2: Generate an offspring population by iterating the following procedures *N* times:

      (1) Select a pair of parent solutions from the current population.

      (2) Generate an offspring from the selected parent solutions by genetic operations.

Step 3: Merge the offspring population and current population. Then select the best *N* solutions from the merged population to construct the next population.

Step 4: If a pre-specified stopping condition is satisfied, terminate the execution of the algorithm. Otherwise, return to Step 2. In the former case, we choose all the non-dominated solutions in the merged population in Step 3 as the final solutions.

Controlled elitist genetic algorithm, a variant of NSGA-II, was proposed by Deb and Goel [20] for controlling the extent of the elite members of the population to maintain the diversity of the population for convergence to an optimal Pareto front. The controlling mechanism is accomplished by allowing only a certain portion of the population to be included in the currently-best-non-dominated solutions. The controlled NSGA-II has a better convergence than the original NSGA-II [20], and since we apply the multi-objective genetic algorithm in different steps of optimisation, we choose to use the controlled NSGA-II in our study rather than the original NSGA-II in order to reduce the computational cost.

**Methodology**

The credit scoring models were implemented using MATLAB 7.5.0. The two proposed methods, which are based on Takagi-Sugero (TS) and Mamdani fuzzy systems are described in the respective order as follows.

*First approach: Takagi-Sugeno-fuzzy-based system*

In the first approach, the fuzzy systems were extracted from the data by a subtractive clustering method and then the resulting fuzzy rules were optimised to increase the accuracy using genetic algorithms. In the last two steps a multi-objective genetic algorithm was applied to preserve the accuracy of the fuzzy model to a given value while enhancing the transparency of the fuzzy model by reducing the customer input and fuzzy sets in the rule base. The steps are outlined below.

**Step 1: Structure and parameter initialisation using subtractive clustering algorithm.** In this step, a fuzzy system of TS type was generated using subtractive clustering method [21] which is an efficient and fast algorithm used for estimating the number of clusters and the location of cluster centres in a set of data. The linear least-square estimation was then used to determine each rule consequent equation. This algorithm has the advantage of describing the TS fuzzy model with few rules [21]. The TS fuzzy model [22] uses fuzzy rules with fuzzy antecedents and functional consequent parts. This model is represented by a series of fuzzy rules of the form:

$$R_k: IF\ x\ is\ A^k\ Then\ y\ is\ f(x) \tag{1}$$

where $R_k$ is the label of the $k^{th}$ fuzzy rule, $f$ represents the output variable $y$, and $A^k$ is the fuzzy set that is defined over input $x$, where $x = (x_1, ..., x_n)$ is the $n$-dimensional pattern vector. $A^k$ is represented by Gaussian membership functions of the form:

$$\mu_{ik}(x_i) = \exp\left(-\frac{(x_i - c_{ik})^2}{2a_{ik}^2}\right) \tag{2}$$

where $c_{ik}$ and $a_{ik}$ are the centre and the width of the Gaussian function respectively.

**Step 2: Structure and parameter optimisation by genetic algorithm.** In the second step, a genetic algorithm was applied to increase the accuracy of the initial fuzzy system by searching for the most suitable value of centre $c_{ik}$ and width $a_{ik}$ of each fuzzy set in the rule base. The fitness function of the genetic algorithm for an individual $S_i$ is given by:

$$Fitns(S_i) = fitn_{acc}(S_i) \tag{3}$$

where $fitn_{acc}(S_i)$ is the fitness function of the accuracy measured by the percentage of correctly classified training patterns. The parameters of membership functions in the antecedents of each fuzzy rule were encoded into a chromosome. Thus, the $i$-th chromosome is a string of the form:

$$S_i = (\underbrace{c_{11}^{(i)}, a_{11}^{(i)}, ..., c_{n1}^{(i)}, a_{n1}^{(i)}}_{premise\ of\ rule\ 1}, ..., \underbrace{c_{1K}^{(i)}, a_{1K}^{(i)}, ..., c_{nk}^{(i)}, a_{nk}^{(i)}}_{premise\ of\ rule\ K}) \tag{4}$$

The first individual of the initial population was generated as a copy of the premise parameters of the initial fuzzy rules generated from step 1. The initialisation from a good population may speed up the convergence of the solution. The remaining individuals were initialised with random values. The best individuals in the population were always selected and kept unchanged in the next generations according to the elitist strategy. The simplest form of crossover, which is the single-point crossover, was adopted. At the end of this step, the highest-accuracy fuzzy model was obtained.

**Step 3: Feature selection using multi-objective genetic algorithm.** The objective of this step is to reduce input dimensions by choosing the relevant subset of features. To achieve this, we need first to keep the accuracy achieved in the previous GA optimisation step as high as possible while choosing the subset which contains the smallest number of features.

The modelling objectives of fuzzy system S in this step can be written as follows:

$$\text{Maximise } f_{acc}(S), \text{ Minimise } f_{input}(S) \tag{5}$$

where $f_{acc}(S)$ is the fuzzy system accuracy measured by the percentage of correctly classified training patterns and $fitn_{input}(S)$ is the total number of selected features of a fuzzy system. To simultaneously achieve these two objectives, controlled elitist genetic algorithm (controlled NSGA II) was applied. The results of this step are Pareto-front solutions that represent a number of fuzzy models with different accuracy numbers of input values. The fuzzy model chosen in this case is based on the need of the user, that is, if the accuracy is more important than the transparency then a fuzzy model with high accuracy and high number of features will be chosen. As stated before, our objective is to enhance the transparency while keeping almost the same accuracy. So a fuzzy model which has accuracy value near to the initial fuzzy model was chosen.

The final result of this step is a fuzzy model with relatively good classification accuracy and relatively fewer number of features.

The following is the design of the chromosome used for feature selection (step3) as well as the genetic operators applied for exploring the search space by testing every possible combination of candidate features and selecting the relevant ones.

Chromosome design

The chromosome $S_i$ which represents the selected features is denoted by a concatenated binary bit string of length $n$ ( $n$ is the total number of features in the data set), where each binary bit denotes whether a given input is selected during the feature selection process. In this implementation, the selected features were set to 1 while the non-selected features were set to 0. Figure 1 shows the structure of an example of one chromosome after the selected feature process in the Australian data set. The total number of features is 14 and the selected inputs which have the value 1 are: 1, 2, 5, 8, 9, 10, 13 and 14.

| F1 | F2 | F3 | F4 | F5 | F6 | F7 | F8 | F9 | F10 | F11 | F12 | F13 | F14 |
|----|----|----|----|----|----|----|----|----|-----|-----|-----|-----|-----|
| 1  | 1  | 0  | 0  | 1  | 0  | 0  | 1  | 1  | 1   | 0   | 0   | 1   | 1   |

**Figure 1.** Chromosome of the genetic algorithm used in the feature selection

Genetic operators

A new population $P$ of chromosomes (fuzzy systems) was generated using the genetic operations: selection, crossover and mutation. To generate a new fuzzy system $S_i$, first a pair of parent fuzzy systems was selected from the current population using tournament selection based on the Pareto ranking and the crowding distance. In order to maintain the diversity in the next population, the best non-dominated solutions were kept down to only 35% (which is the default value defined by the algorithm [23]) of the population. In addition, the crowding measure was used to calculate the crowding distance for each individual on a non-dominated front. After the selection step, the uniform crossover and uniform mutation with a range of 0.01 were applied. These genetic operations were applied for fuzzy sets selection step (next step) and also for fuzzy sets replacement process in the second approach because they involve similar problem.

**Step 4: Fuzzy sets selection using multi-objective genetic algorithm.** The final step in the optimisation process is the removal of the fuzzy sets whose effect on the fuzzy system is not important. In this case, we have two objectives; the first one is to maximise the classification accuracy of the fuzzy system while minimising the number of fuzzy sets in the fuzzy system is the second objective. The objectives of the fuzzy system S are written as follows:

Maximise $f_{acc}(S)$ , Minimise $f_{sets}(S)$ (6)

where $f_{acc}(S)$ is the accuracy of the fuzzy system measured by the percentage of correctly classified training patterns and $f_{sets}(S)$ is the number of selected fuzzy sets. To accomplish these two objectives, controlled elitist genetic algorithm was applied. Among the Pareto-front solutions

that represent fuzzy models with different accuracy numbers of fuzzy set trade-off, a fuzzy system with accuracy almost equal to the initial fuzzy system was chosen. After this final stage, a fuzzy model with relatively good classification accuracy and less fuzzy sets was obtained.

The following is a description of the chromosome design used for fuzzy set selection (step 4). In this step, we adopt the same genetic operators as in step 3.

Chromosome design

The chromosome $S_i$ which represents the selected antecedents is denoted by a concatenated binary bit string of length $L_i = n' \times K$, where $n'$ and $K$ are the number of inputs after feature selection phase and the number of rules in the fuzzy system respectively. Each binary bit in the string denotes whether a given fuzzy antecedent is selected. In this case, the selected antecedents were set to 1 and non-selected antecedents were set to 0. Figure 2 shows the structure of the chromosome used in the antecedent fuzzy sets selection phase.
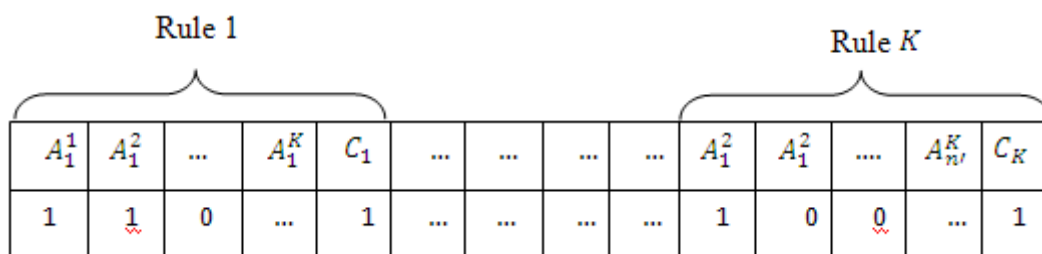


**Figure 2.** Chromosome of the genetic algorithm used in the antecedent fuzzy set selection phase

*Second approach: Mamdani-fuzzy-based system*

In the second approach, we used Mamdani fuzzy system [24] in place of Takagi-Sugeno one. First, the initial Mamdani fuzzy system was generated using fuzzy C-means clustering (FCM) method [25]. The generated Mamdani fuzzy rules are written as:

$$R_k: IF\ x\ is\ A^k\ Then\ y\ is\ B^k \qquad (7)$$

where $R_k$ is the label of the $k^{th}$ fuzzy rule, $A^k$ is the fuzzy set defined over the input $x$ where $x = (x_1, ..., x_n)$ is the $n$-dimensional pattern vector while $B^k$ is a fuzzy set defined over the output variable $y$. All the fuzzy sets in the rule base are represented by Gaussian function.

The next step is to replace the fuzzy sets of the generated fuzzy system by new fuzzy sets. The reason behind this replacement is that the fuzzy sets resulting from clustering or learning method are usually not interpretable [10]. On the other hand, the new predefined fuzzy sets have clear linguistic interpretations such as low, average and high. The linguistic values of each attribute $x_i$ have to be defined before starting the replacement process. In our case, we use five linguistic values: very low, low, average, high and very high, and each of the linguistic values is defined within a specific range of values. Figure 3 shows an example of five linguistic values for the credit amount attribute in German credit data. These new fuzzy sets replace the existing ones

of the credit amount attribute in the fuzzy-rule-based system. This idea is similar to that applied by Ishibuchi et al [26].

In replacing the existing fuzzy sets by the new ones the following must be considered:
- The replacement of an existing fuzzy set $A^{ik}$ of the rule $R_k$ and the $x_i$ attributes by $A^{,ik}$ where $A^{,ik}$ is one of the linguistic values defined over the $x_i$ attribute. (For example, $A^{,ik}$ could be either low, average or high.)
- The replacement procedure has to improve the classification accuracy of the fuzzy system.
- In addition to the five linguistic values, 'don't care' is another linguistic value and it refers to unimportant fuzzy set that can be deleted without effecting the fuzzy system performance.

In the first subsection below more explanation on the problem of replacing the existing fuzzy sets with linguistic values is given and the proposed solution is described. In the second subsection, a description of the chromosome design used for Mamdani-based fuzzy system is given.
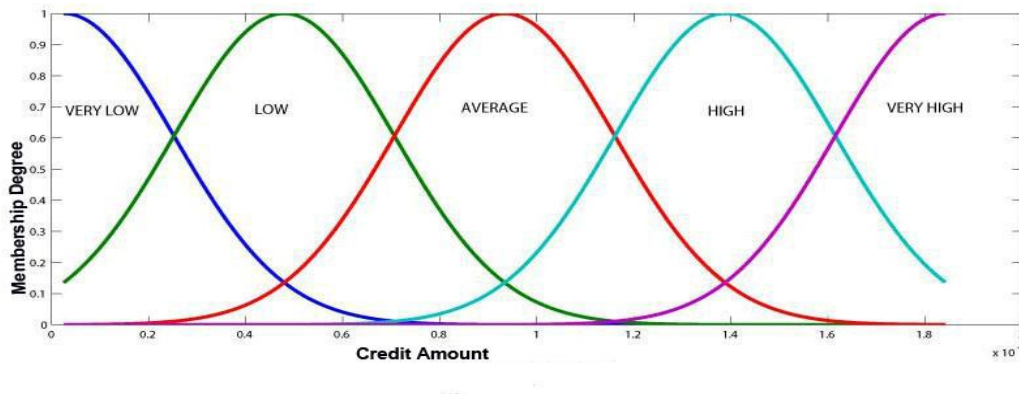


**Figure 3.** Linguistic fuzzy sets of the credit amount attribute – German data set

Problem formulation

Let $K_i$ be the number of linguistic values in each attribute. So, for each attribute $x_i$ , we have $K_i$ possible antecedent fuzzy sets. In addition, 'don't care' is considered as another fuzzy set. In this case, we have $(K_i + 1)$ possible cases and each antecedent fuzzy set $A^{ik}$ in the fuzzy rules is selected from the given $K_i$ linguistic values and 'don't care'. The total number of possible combinations of the antecedent linguistic values in the fuzzy rules is $(K_{11} + 1) * (K_{22} + 1) * ... * (K_{in} + 1)$, where $K_{in}$ is the number of linguistic values of $x_i$ and $n$ is the number of antecedent fuzzy sets in the fuzzy rules.

The task now is to search for the best combination of these antecedent linguistic values that achieves the two objectives, namely maximising the classification accuracy and maximising the transparency by increasing the number of 'don't care' fuzzy sets in the rule base. These two objectives of the fuzzy system S can be written as:

$$\text{Maximise } f_{acc}(s), \text{ Maximise } f_{transp}(s) \qquad (8)$$

where $f_{acc}$(s) is the classification accuracy of the fuzzy system measured by the percentage of correctly classified training patterns and $f_{transp}$ is the transparency measured by the number of 'don't care' fuzzy sets in the rule base. To solve this combinatorial problem with these two objectives, a controlled elitist genetic algorithm is applied. Since 'don't care' conditions can be omitted, fuzzy rules with many 'don't care' conditions are written as short fuzzy rules.

Chromosome design

The chromosome $S_i$ is coded as follows:

$$S_i = (\ \underbrace{A_1^1\ A_2^1\ A_3^1\quad \ldots\quad A_n^1\ C_1}_{Premise\ and\ consequent\ of\ rule\ 1}\ ,\ldots,\ \underbrace{A_1^K\ A_2^K\ A_3^K\quad \ldots\quad A_n^K\ C_K}_{Premise\ and\ consequent\ of\ rule\ K}\ ) \tag{9}$$

where $n$ and $K$ denote the number of features and fuzzy rules respectively. $A_i^j$ is the linguistic antecedent value and $C_K$ is the consequent class. The length of the chromosome is $(n+1) \times K$. We used six linguistic values: very low, low, average, high, very high and 'don't care'. Each of these linguistic values is defined by a number. In our case, we set the values 0, 1, 2, 3, 4 and 5 to denote 'don't care', very low, low, average, high and very high respectively. For the consequent class, we set 0 and 1 for negative and positive class respectively. In this case, each antecedent condition $A_i^j \in \{0,1,2,3,4,5\}$ and the consequent class $C_K \in \{0,1\}$.

The following is an example to further explain this idea. Assume that we generate a fuzzy system with 3 inputs and 2 rules in the clustering step and then we get the string 01204521 as one of the best Pareto solutions at the end of the multi-objective optimisation process. Figure 4 shows the decoding process of the 01204521 string. Since we have 3 inputs and one output, the length of string encoding one rule is four. Decoding process of the previous string results in the following rules:

Rule1: If input 1 is 'don't care', input 2 is very low and input 3 is low, then outcome is negative.
Rule2: If input 1 is high, input 2 is very high and input 3 is low, then outcome is positive.
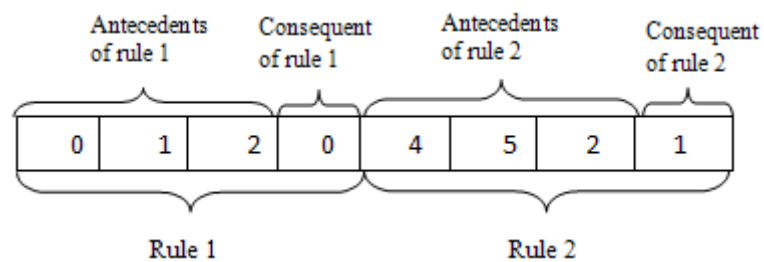


**Figure 4.** Chromosome coding with 3 inputs and 2 rules

**Data Sets**

We used two data sets, namely Germany credit data set [13] and Australian credit data set [14]. Both data sets are made publicly available to all users from the UCI Repository of Machine Learning Databases [13] and are mostly used to compare the performance of various classification models.

To eliminate the skewness and bias in the training and test samples, a common validation method called random sub-sampling validation was applied [27]. Using this method, the whole of the two data sets (German and Australian) were randomly divided into two parts: one for training and the other for testing. As commonly done in similar studies [27], 70% of the data were used for training purpose and 30% for testing the performance of the model. This process was repeated five times to create five pairs of training and test samples for each data set. To calculate the fuzzy classification accuracy for example, five fuzzy systems were built using the training data sets of the five data partitions (S1, S2, S3, S4 and S5) and the accuracy result (validation result) was averaged over the classification accuracy of the corresponding five test sets. The advantage of this technique over *k*-fold and leave-one-cross validation methods is that the proportion of the training/testing split is independent of the number of iterations (folds). Table 1 presents the features of the data sets used in this study. (See Tables 9 and 10 in Appendix for more details on the data attributes.)

**Table 1.** German and Australian data sets

|  | No. of inputs | Training set size | Testing set size | Data set size |
|---|---|---|---|---|
| German data | 20 | 700 | 300 | 1000 |
| Australian data | 14 | 383 | 207 | 690 |

## Results and Discussion

*The first approach: TS-fuzzy-based system*

The results of the first step are summarised in Table 2. Since the transparency of the credit scoring model is one of the two modelling objectives, a compact fuzzy system with 3 fuzzy rules and relatively good accuracy was generated in the first step for both data sets.

**Table 2.** Classification accuracy of initial fuzzy systems with 3 rules using five randomly generated samples

|  |  | S1 | S2 | S3 | S4 | S5 | Average |
|---|---|---|---|---|---|---|---|
| Australian data set | Training accuracy (%) | 88.61 | 88.82 | 89.03 | 88.82 | 87.58 | 88.57 |
|  | Testing accuracy (%) | 89.37 | 86.96 | 86.96 | 86.96 | 86.47 | 87.34 |
| German data set | Training accuracy (%) | 75.57 | 75.71 | 77.86 | 76.29 | 77.43 | 76.57 |
|  | Testing accuracy (%) | 75.00 | 69.67 | 74.67 | 73.33 | 71.33 | 72.80 |

In the second step, a genetic algorithm was applied to improve the performance of the fuzzy system generated from the first step. The disadvantage of the genetic algorithm is its computational cost, but starting from a good point (initial fuzzy model) has speeded up the convergence of the genetic algorithm. For example, the genetic algorithm in the case of the Australian data-S1 took around 15 epochs to converge while it took around 16 epochs in the case

of the German data-S1. This reveals the complementary functioning between fuzzy clustering and genetic algorithm. The results of this step are shown in Table 3. As the table indicates, there is an increase in the prediction accuracy for both data sets. For the Australian data set, the classification accuracy for testing data increases from 87.34% to 88.89% while that for the German data set has a relatively significant increase from 72.80% to 77.07%. The difference in the enhancement may be due to the difference in the degree of the complexity in the two data sets. The German data is more complicated than the Australian one as the former has 20 inputs while the latter has 14 inputs. Thus, the genetic algorithm seems to be more efficient in dealing with this complexity than the subtractive clustering algorithm.

In the third step, a multi-objective genetic algorithm was applied to select the most relevant inputs in both data sets. Table 4 summarises the results of this step. The names and characteristics of the selected inputs for the two data sets are listed in Tables 9-10 (Appendix). For example, as Table 4 shows, the selected inputs for the German data-S1 are: 1, 2, 3, 5, 8, 12, 14 and 19 and the corresponding attributes are the following: (1) Status of existing checking account, (2) Duration in month, (3) Credit history, (5) Credit amount, (8) Installment rate, (12) Property, (14) Other installment plans, and (19) Telephone. The Australian data attributes have been changed to meaningless symbols to protect the confidentiality of the data. As can be seen from Table 4, there is an improvement in the transparency of the fuzzy systems for both data sets represented by the decrease in the number of the inputs while there is a slight decrease in the accuracy of the fuzzy systems. Particularly, in the Australian data set the accuracy is kept almost the same. Its average number of inputs of fuzzy models is 5.8 while it is 8.6 for the German data set. The new fuzzy system of the Australian data-S1 has only 6 inputs and the other inputs have been deleted without affecting the prediction accuracy of the model. Hence, there is a complexity (in this case unnecessary attributes) that should be removed without decreasing the model performance and that has no relation with the accuracy-transparency trade-off. Furthermore, there are some cases where removing some inputs may increase the accuracy (like in the case of Australian data-S5). On the other hand, results from the German data case reveal accuracy-transparency trade-off. For example, the performance of the fuzzy system of sample 3 that contains 11 inputs is better than that of the other fuzzy models with fewer numbers of inputs.

**Table 3.** Classification accuracy of fuzzy systems after applying genetic algorithm on initial fuzzy systems

| | | S1 | S2 | S3 | S4 | S5 | Average |
|---|---|---|---|---|---|---|---|
| Australian data set | Training accuracy (%) | 88.82 | 88.00 | 90.06 | 89.65 | 89.44 | 89.19 |
| | Testing accuracy (%) | 89.86 | 87.92 | 88.41 | 89.37 | 88.89 | 88.89 |
| German data set | Training accuracy (%) | 78.00 | 77.71 | 78.86 | 77.43 | 78.29 | 78.06 |
| | Testing accuracy (%) | 78.00 | 74.67 | 77.67 | 78.00 | 77.00 | 77.07 |

In the final step, the unnecessary fuzzy sets were removed. Table 5 shows the results of this step for both data sets. For the Australian data set, the average number of fuzzy sets per rule decreases from 5.8 to 3 while the prediction accuracy is maintained as it was before transparency

optimisation. In the case of the German data set, there is an improvement in the transparency of the fuzzy system from 8.6 to 4.9 fuzzy sets per rule with a slight decrease in the prediction accuracy. The fuzzy rules resulting from the last step of the Australian data-S1 are depicted in Figures 5-6. As Table 5 shows, 8 fuzzy sets are removed with only a slight decrease in the classification accuracy from 89.86% to 89.37%.

**Table 4.** Classification accuracy of fuzzy systems and their corresponding number of selected inputs after applying feature selection procedure

| | | S1 | S2 | S3 | S4 | S5 | Average |
|---|---|---|---|---|---|---|---|
| Australian data set | No.of selected inputs | 6 | 5 | 6 | 7 | 5 | **5.8** |
| | Inputs selected | 1, 4, 8, 9, 11, 13 | 1, 3, 4, 10, 11 | 1, 4, 7, 9, 11, 13 | 1, 3, 6, 8, 9, 10, 11 | 2, 7, 8, 11, 14 | |
| | Training accuracy (%) | 87.58 | 87.58 | 87.37 | 88.61 | 88.41 | **87.91** |
| | Testing accuracy (%) | 89.86 | 87.92 | 87.92 | 88.41 | 89.37 | **88.70** |
| German data set | No.of selected inputs | 8 | 9 | 11 | 6 | 9 | **8.6** |
| | Inputs selected | 1, 2, 3, 5, 8, 12, 14, 19 | 1, 2, 5, 8, 9, 10, 11, 15, 18 | 1, 3, 6, 7, 9, 10, 13, 14, 15, 18, 19 | 1, 5, 7, 10, 16, 17 | 1, 2, 4, 6, 9, 10, 13, 17, 19 | |
| | Training accuracy (%) | 75.00 | 75.42 | 76.57 | 75.43 | 75.57 | **75.60** |
| | Testing accuracy (%) | 75.33 | 73.00 | 77.00 | 75.00 | 75.67 | **75.20** |

**Table 5.** Classification accuracy of fuzzy systems and their corresponding number of selected fuzzy sets after applying fuzzy set selection procedure

| | | S1 | S2 | S3 | S4 | S5 | Average |
|---|---|---|---|---|---|---|---|
| Australian data set | No.of selected fuzzy sets (sets/rule) | 10 | 12 | 8 | 8 | 7 | **3** |
| | Training accuracy (%) | 88.00 | 87.58 | 88.20 | 88.82 | 88.41 | **88.20** |
| | Testing accuracy (%) | 89.37 | 88.41 | 87.92 | 88.89 | 89.41 | **88.60** |
| German data set | No.of selected fuzzy sets (sets/rule) | 12 | 12 | 14 | 17 | 18 | **4.9** |
| | Training accuracy (%) | 75.29 | 74.57 | 75.43 | 75.71 | 76.57 | **75.51** |
| | Testing accuracy (%) | 75.67 | 73.00 | 75.33 | 75.00 | 76.00 | **75** |

*(1) IF Input4 is Gaussian(0.6656 2.042) AND IF In-put8 is Gaussian(0.3767 0.03808) AND IF Input11 is Gaussian (0.3998 0.1251) AND IF Input13 is Gaussian(305.7 199.7)Then The customer is GOOD*

*(2) IF Input9 is Gaussian(0.3874 0.9673) AND IF Input11 is Gaussian(0.4078 0.9664) Then The customer is BAD*

*(3) IF Input1 is Gaussian(0.3679 0.2458) AND IF Input4 is Gaussian(0.6675 2.005) AND IF Input9 is Gaussian (0.3693 0.8825) AND IF Input13 is Gaussian(312.1 0.06008) Then The customer is BAD*

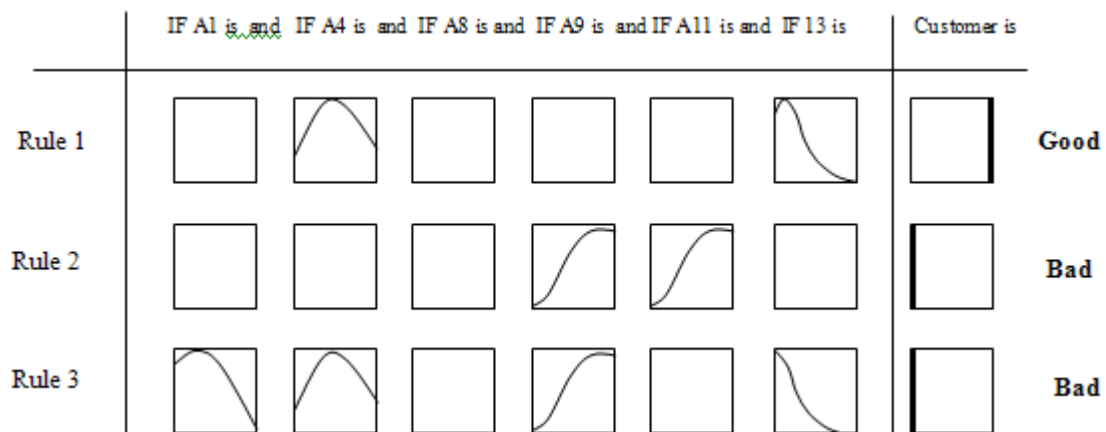**Figure 5.** Approximate TS fuzzy rules generated from Australian data-S1



**Figure 6.** Fuzzy rules after the fuzzy set selection step (Australian data-S1)

**Evaluation of the transparency of the antecedent fuzzy sets.** The fuzzy IF-THEN rules depicted in Figure 5 are approximate fuzzy rules and they do not give an accurate description of the antecedent of the fuzzy rules and therefore this fuzzy system is not considered transparent despite its compactness, as it does not satisfy one of the comprehensibility measures which is the linguistic representation of the produced fuzzy rules. However, some useful information like the attributes that influence or contribute to the system decision can be extracted and some of the antecedent fuzzy sets can be easily understood using the graphical plot of fuzzy rules and even transformed into linguistic values. For example, the 11th attribute of the Australian data set is categorical and has two values; the first one is 0 and the second is 1. As Figure 5 shows, this attribute is only included in the first and second rules that represent the good and bad classes respectively. It can be seen that the value of the 11th attribute is 0 in the first rule while it takes the value 1 in the second rule and it can be concluded that the 11th attribute of the Australian data set has an effect in the discrimination between the good and bad customers. Furthermore, linguistic values can be assigned to both of the two fuzzy sets. Therefore, rather than saying IF Input 11 is Gaussian (0.3998 0.1251) in the first rule and IF Input 11 is Gaussian (0.4078 0.9664) in the second rule, we can say IF Input 11 belongs to category 1 (value 0) and IF Input 11

belongs to category 2 (value 1) in the first and the second rule respectively. Another issue for the continuous attribute is when there is too much overlapping between the antecedent fuzzy sets. This problem prevents the readability of the fuzzy sets and causes a lack of comprehensibility of the fuzzy system. To overcome this problem, a similarity-driven procedure [28] can be applied to merge similar antecedent fuzzy sets into a given attribute.

**Evaluation of the accuracy of the TS fuzzy system.**   To evaluate the performance of our approach, the credit scoring model developed in this study is compared with the other benchmark methods [1, 29-30] applied on the same data sets. These methods are usually used to test the classification accuracy of the new algorithms applied for credit scoring models. More details about the main characteristics of these methods and their application in the credit scoring systems are described by Lahsasna et al [3]. As Table 6 shows, the first approach applied in this study (TS fuzzy system) compares favourably with the other methods such as genetic programming (GP), artificial neural networks (ANNs), radial basis function (RBF) and genetic algorithms-support vector machines ('GA+SVM') hybrid method while it is superior to some methods such as classification and regression tree (CART), rough sets, and the popular decision tree algorithm C4.5. Even though the machine learning methods are accurate classification methods, the lack of transparency of these methods is a major drawback especially when the end user needs to get some information about the credit system. Unlike these black-box methods like ANNs, SVM and genetic algorithm, TS-fuzzy-based method gives some useful information (such as defining the customer's attributes) that influences the system decision and the approximate values of these attributes.

**Table 6.**   Classification accuracy of 'GA+SVM', GP, CART, C4.5, rough sets, ANNs, RBF and TS-fuzzy-based system

| Author | Method | Classification accuracy for Australian data (%) | Classification accuracy for German data (%) |
|--------|--------|:-----:|:-----:|
| [29] | GA+SVM | 86.9 | 77.92 |
| [30] | GP | 88.27 | 77.34 |
| [30] | CART | 85.81 | 70.59 |
| [30] | C4.5 | 87.06 | 73.17 |
| [30] | Rough sets | 83.72 | 74.57 |
| [30] | ANNs | 87.93 | 75.51 |
| [1] | RBF | 87.78 | 75.63 |
| This paper | TS-fuzzy-based system (accuracy only) | 88.89 | 77.07 |

In the case where the end user is only interested in prediction (i.e. getting the best classification accuracy), the use of TS-fuzzy-based system which results from step 2 (structure and parameter optimisation by genetic algorithm) is recommended as it is more accurate than the

TS-fuzzy-based system resulting from the final step where both the accuracy and the transparency have been considered. Alternatively, other methods such as GP, ANNs, RBF and 'GA+SVM' can be used in such case.

*The second approach: Mamdani-fuzzy-based system*

The results obtained at the end of this step are summarised in Table 7, which shows the degree of performance and the level of compactness of each fuzzy system using the rate of correctly classified testing patterns and the number of antecedents per rule respectively. The number of antecedent fuzzy sets per rule has been reduced from 20 and 14 to 6.7 and 3.8 fuzzy sets per rule for the German and Australian data sets respectively. So this method gives relatively good results in enhancing the compactness of the initial fuzzy system resulting from the clustering step.

**Table 7.** Accuracy and transparency results for Mamdani-based fuzzy system

| | | S1 | S2 | S3 | S4 | S5 | Average |
|---|---|---|---|---|---|---|---|
| German data set | Training accuracy (%) | 78.28 | 78.57 | 76.57 | 77.14 | 78.86 | **77.88** |
| | Testing accuracy (%) | 74.33 | 71 | 73.33 | 72.33 | 72 | **72.60** |
| | Total no.of fuzzy sets in fuzzy system | 160 | 160 | 160 | 160 | 160 | **160** |
| | Total no.of selected fuzzy sets in fuzzy system | 56 | 61 | 53 | 49 | 50 | **53.8** |
| | Average no.of selected fuzzy sets per rule | 7 | 7.6 | 6.6 | 6.1 | 6.25 | **6.71** |
| Australian data set | Training accuracy (%) | 86.75 | 87.78 | 90.26 | 86.75 | 88 | **87.91** |
| | Testing accuracy (%) | 88.88 | 86 | 86 | 86 | 84.05 | **86.19** |
| | Total no.of fuzzy sets in fuzzy system | 98 | 98 | 98 | 98 | 98 | **98** |
| | Total no.of selected fuzzy sets in fuzzy system | 19 | 26 | 25 | 36 | 29 | **27** |
| | Average no.of selected fuzzy sets per rule | 2.71 | 3.71 | 3.57 | 5.14 | 3.6 | **3.8** |

**Evaluating the antecedent fuzzy sets comprehensibility.** The antecedent fuzzy sets of this system become well defined and distinguishable. In Figure 3, the antecedent fuzzy sets of the attribute credit amount are plotted. This attribute may be assigned five well defined linguistic values: very low, low, average, high and very high, and every linguistic fuzzy set has a specific range of values. The fuzzy rules are humanly understandable because of the natural language used. Hence, the descriptive Mamdani fuzzy rules generated have the capacity to represent the knowledge characterising the relations between the customer features and his creditworthiness in a series of linguistic fuzzy rules, thus rendering the decision process of the system understandable

and helping the manager in the financial institution to make useful financial analysis and then make the right decisions.

**Comparison between TS and Mamdani fuzzy systems.** Table 8 shows accuracy and transparency results for German and Australian data sets using TS and Mamdani fuzzy systems. Transparency results are shown using transparency 1, which indicates the number of rules, and transparency 2, which defines the number of fuzzy sets per rule. In addition, Figure 7 shows Mamdani fuzzy rules generated from Australian data set while Figure 8 shows the fuzzy rules extracted from German data set. The fuzzy system uses IF-THEN rules with linguistic values. By comparing this rule set with the rules extracted from TS fuzzy approach using the performance and comprehensibility criteria, the following results are noted.

*Performance* : TS fuzzy system performance is 88.60% and 75% for the Australian and German data sets respectively while for Mamdani fuzzy system it is 86.19% and 72.60% for the same data sets. These results indicate that TS fuzzy system is more powerful than Mamdani fuzzy system and thus the former is the better choice for predicting the customer's creditworthiness.

*Comprehensibility* : Compared to Mamdani fuzzy system TS fuzzy system is more compact as it uses only 3 rules for both data sets while the former uses 7 rules for Australian data and 8 rules for German data. Furthermore, TS fuzzy system generally uses slightly smaller number of antecedent fuzzy sets per rule than Mamdani system for both data sets. Despite these strong points of the TS fuzzy system, however, the Mamdani fuzzy system has a major advantage in the capacity to represent the fuzzy rules in an intuitive way using linguistic fuzzy rules. This capacity represents the true level of comprehensibility. The approximate fuzzy rules of the TS fuzzy system do not give a clear idea about the underlying relation between the customer features and their creditworthiness or generally between the input and the output of the data.

**Table 8.** Accuracy and transparency results for TS- and Mamdani-based fuzzy systems

| Data | TS-fuzzy-based system | | | Mamdani-fuzzy-based system | | |
|---|---|---|---|---|---|---|
| | **Accuracy** | **Transparency 1** | **Transparency 2** | **Accuracy** | **Transparency 1** | **Transparency 2** |
| German | 75% | 3 rules | 4.9 sets/rule | 72.60% | 8 rules | 3.8 sets/rule |
| Australian | 88.60% | 3 rules | 3 sets/rule | 86.19% | 7 rules | 6.71 sets/rule |

Despite the enhancement of the comprehensibility of approximate of TS fuzzy sets using similarity-driven method, the problem is not definitely resolved especially when there is a high number of fuzzy sets in the same attribute. The Mamdani fuzzy system is therefore the best choice for data analysis and knowledge discovery from the data set. Thus, the choice of the system type is based on the needs of the user as to whether a high accuracy prediction or a high comprehensibility system is needed. For generating a completely transparent credit scoring model, Mamdani fuzzy system should be chosen. This makes the credit scoring model easier to

(1)  (IF A4 is p(=1)) and (IF A9 is t(=0))(Then  customer is bad)

(2)  (IF A5 is j(=5)) and (IF A9 is f(=1)) and (IF A12 is s(=1))(Then  customer is bad)

(3)  (IF A5 is j(=5)) and (IF A8 is f(=1))(Then customer is good)

(4)  (IF A5 is aa(=6)) and (IF A9 is t(=0)) and (IF A10 is very low) and (IF A11 is f(=1))(Then customer is bad)

(5)  (IF A5 is c(=8)) and (IF A10 is very low) and (IF A14 is average )(Then  customer is bad)

(6)  (IF A4 is p(=1)) and (IF A8 is f(=1)) and (IF A12 is p(=3)) and (IF A13 is very low)(Then customer is good)

(7)  (IF A13 is very high)(Then customer is bad)

**Figure 7.**  Descriptive Mamdani fuzzy  rules  generated from  Australian data-S1

**(1)** (IF Credit amount is very low) and (IF unknown/no savings account) and (IF Other debtors/guarantors: guarantor) and (IF  Age is very low) and (IF  Other installment plans: bank) and (IF Housing :for free)(Then  customer is bad)

**(2)**  (IF existing credits paid back ) and (IF Purpose is radio/television) and (IF savings < 100  DM) and (IF employment since :unemployed) and (IF male: single) and (IF Property is life insurance) and (IF Housing :rent) and (IF Number  of  credits=1)(Then customer is good)

**(3)** (IF Purpose is car (used)) and (IF unknown/no savings account)  and (IF  Other debtors/guarantors: none) and (IF Property car or  other) and (IF Age is very low)(Then customer is good)

**(4)** (IF Check Account ≥ 200   DM ) and (IF  Purpose is vacation)  and (IF  Credit amount is low)  and  (IF male: married/widowed) and (IF Other debtors/guarantors :none) and (IF Residence=4 years) and (IF Property is life insurance) and (IF Age is very low) and (IF Housing :rent)(Then customer is bad)

**(5)** (IF Check Account < 0  DM) and (IF critical account)and (IF Purpose is education) and (IF savings≥ 1000 DM) and (IF  4  ≤ employment since <7 years) and (IF Other debtors/guarantors:none)  and (IF  Other installment  plans:bank)  and  (IF  Housing :own)  and  (IF Number  of  credits=4)(Then customer is bad)

**(6)**  (IF  Duration  is  too  short)  and  (IF  critical  account) and (IF 4≤employment since<7 years) and (IF female:divorced/separated/married) and (IF Other debtors/guarantors:none) and (IF Residence=4 years) and (IF Housing :rent) and (IF Number of  credits=1) and (IF Job is unskilled - resident)(Then  customer is good)

**(7)** (IF Check Account ≥ 200  DM ) and (IF Duration is very long) and (IF  500   ≤ savings <1000 DM ) and (IF Other debtors/guarantors:none) and (IF  Age  is  very low)(Then customer is bad)

**(8)**  (IF  delay  in  paying  off  in  the  past)  and  (IF male:married/widowed) and (IF Property car or  other) and (IF  Age is high) and (IF  Number of  people=1)(Then  customer is good)

**Figure 8.**  Descriptive Mamdani fuzzy  rules  generated from  German  data-S1

understand, for example the reason behind some decisions such as rejecting a credit application. In such a case, it is reasonable to trade some accuracy for extra transparency and better readability of the credit scoring model. The multi-objective genetic algorithm applied in this study can achieve a maximum trade-off between the accuracy and transparency. Hence, an adequate credit scoring system can be chosen based on the needs of the user, for example in the case where the end user wants only to conduct a data analysis to find out about the main customer attributes that influence the discrimination between the good and bad customers. In this case, the transparency which is measured by the number of selected inputs is more important than the classification accuracy and the recommended choice for him/her is to select an accuracy-transparency level where the accuracy value is acceptable while the transparency value is as high as possible (i.e. select the minimum number of inputs). The acceptable level of accuracy is the minimum level required to have a reliable data analysis while a very high level of transparency allows for better understanding of the key factors that influence the classification process. In another case in which the end user is only interested in the outcome without paying attention to the interpretation of the results, the fuzzy credit system with the highest classification accuracy value is chosen, irrespective of the number of selected inputs. The first approach (TS fuzzy system) is suitable for getting the above-mentioned choices. To further the investigation and the analysis on the relation between the attributes and the outcome, the end user needs to see the variation in the outcome when the values of certain attributes change. In such a case, the values of the fuzzy sets have to be well-defined and distinguishable so that each of the fuzzy sets can be defined using a linguistic value. The linguistic values such as low, average and high are natural and humanly understandable values and can be used as labels for the fuzzy sets to construct the fuzzy system. The second approach (Mamdani fuzzy system) is suitable for this kind of data analysis where the end user is interested in knowing not only the important attributes that contribute to the outcome but also the details on the relation between the attributes and the outcome.

## Conclusions

In this paper, the transparency and accuracy of credit scoring model have been investigated using two different fuzzy model types, namely Takagi-Sugeno (TS) and Mamdani. The following conclusions have been drawn from this study.

TS fuzzy system is highly accurate and computationally efficient although lacking in transparency while Mamdani fuzzy system is highly transparent, intuitive, well suited to human input and relatively accurate. Therefore, TS fuzzy system is apparently better in predicting the customer's creditworthiness while the latter is better in data analysis and knowledge discovery. The transparency of the fuzzy systems resulting from clustering techniques is often lost during the learning of parameters and can be evaluated in two levels. The first and most important level of transparency is the capacity to represent the knowledge characterising the relations between the customers' features and creditworthiness in a natural manner, e.g. as a series of linguistic fuzzy rules. The second level is the degree of complexity of the fuzzy system which can be measured by the number of fuzzy rules in the fuzzy system, the number of input variables for each rule, and the

number of fuzzy sets per variable. This study illustrates a classical trade-off between accuracy and transparency, and the power of multi-objective learning to find an adequate trade-off between them so the user can choose between different levels of accuracy-transparency based on the end user's needs. This technique can also remove unnecessary complexity that may result from extracting the rules from the data set without affecting the performance of the fuzzy system. Therefore, it can be used in the pre-processing stage of a high dimensional pattern modelling as a feature selection method to reduce the number of inputs in the model. In this case the user can, based on his need, choose between different levels of number of inputs/accuracy of the model.

## Acknowledgement

## References

1. D. West, "Neural network credit scoring models", *Comput. Operat. Res.*, **2000**, *27*, 1131-1152.

2. L. C. Thomas, "A survey of credit and behavioral scoring: Forecasting financial risks of lending to customers", *Int. J. Forecast.*, **2000**, *16*, 149-172.

3. A. Lahsasna, R. N. Ainon and T. Y. Wah, "Credit scoring models using soft computing methods: A survey", *Int. Arab J. Inform. Technol.*, **2010**, *7,* 115-123.

4. H. L. Jensen, "Using neural networks for credit scoring", *Managerial Finance*, **1992**, *18*, 15-26.

5. V. S. Desai, J. N. Crook and G. A. Overstreet, "A comparison of neural networks and linear scoring models in the credit union environment", *Eur. J. Operat. Res.*, **1996**, *95*, 24-37.

6. B. Baesens, R. Setiono, C. Mues and J. Vanthienen, "Using neural network rule extraction and decision tables for credit-risk evaluation", *Manage. Sci.*, **2003**, *49*, 312-329.

7. S. Piramuthu, "Financial credit-risk evaluation with neural and neuro-fuzzy systems", *Eur. J. Operat. Res.*, **1999**, *112*, 310-321.

8. R. Malhotra and D. K Malhotra, "Differentiating between good credits and bad credits using neuro-fuzzy system", *Eur. J. Operat. Res.*, **2002**, *136*, 190-211.

9. F. Hoffmann, B. Baesens, J. Martens, F. Fput and J. Vanthienen, "Comparing a genetic fuzzy and a neuro-fuzzy classifier for credit scoring", *Int. J. Intell. Syst.*, **2002**, *17*, 1067-1083.

10. D. Nauck, U. Nauck and R. Kruse, "Generating classification rules with neuro-fuzzy system NEFCLASS", Proceedings of the Biennial Conference of the North American Fuzzy Information Processing Society, **1996**, Berkeley, CA, USA, pp.466-470.

11. F. Herrera, "Genetic fuzzy systems: taxonomy, current research trends and prospects", *Evol. Intell.*, **2008**, *1*, 27-46.

12. H. Ishibuchi, "Multiobjective genetic fuzzy systems: review and future research directions", Proceedings of the 2007 IEEE International Conference on Fuzzy Systems, **2007**, London, UK, pp. 913–918.

13. A. Asuncion and D. J. Newman, "UCI Machine Learning Repository", School of Information and Computer Science, University of California, Irvine, CA, http://www.ics.uci.edu/~mlearn/ MLRepository.html. (Retrieved: May 15, **2007**)

14. R. Quinlan, "C4.5: Programs for Machine Learning", Morgan Kaufmann, San Fransisco, **1992**.

15. A. Konak, D. W. Coit and A. E. Smith, "Multi-objective optimization using genetic algorithms: A tutorial", *Reliab. Eng. Syst. Saf.*, **2006**, *91*, 992-1007.

16. C. A. C. Coello, "A comprehensive survey of evolutionary-based multi-objective optimization techniques", *Knowl. Inform. Syst.*, **1999**, *1*, 269-308.

17. D. A. Van Veldhuizen and G. B. Lamont, "Multi-objective evolutionary algorithms: analyzing the state-of-the-art", *Evol. Comput.*, **2000**, *8*, 125-147.

18. K. Deb, A. Pratap, S. Agarwal and T. Meyarivan, "A fast and elitist multi-objective genetic algorithm : NSGA-II," *IEEE Trans. Evol. Comput.*, **2002**, *6*, 182-197.

19. N. Srinivas and K. Deb, 'Multi-objective optimization using non-dominated sorting in genetic algorithms", *Evol. Comput.*, **1994**, *2*, 221-248.

20. K. Deb and T. Goel, "Controlled elitist non-dominated sorting genetic algorithms for better convergence", Proceedings of the 1st International Conference on Evolutionary Multi-Criterion Optimization, **2001**, Berlin, Germany, pp. 67-81.

21. S. L. Chiu, "Fuzzy model identification based on cluster estimation", *J. Intell. Fuzzy Syst.*, **1994**, *2*, 267-278.

22. T. Takagi and M. Sugeno, "Fuzzy identification of systems and its application to modeling and control", *IEEE Trans. Syst. Man Cybern.*, **1985**, *15*, 116-132.

23. The MathWorks, "Matlab Global Optimization Toolbox 3, User's Guide", The MathWorks Inc., Natick, MA, **2010**.

24. E. H. Mamdani, "Applications of fuzzy logic to approximate reasoning using linguistic synthesis", *IEEE Trans. Comput.*, **1977**, *26*, 1182-1191.

25. J. C. Bezdek, "Pattern Recognition with Fuzzy Objective Function Algoritms", Plenum Press, New York, **1981**.

26. H. Ishibuchi, T. Nakashima and T. Murata, "Three-objective genetics-based machine learning for linguistic rule extraction", *Inform. Sci.*, **2001**, *136*, 109-133.

27. R. Kohavi, "A study of cross-validation and bootstrap for accuracy estimation and model selection", Proceedings of International Joint Conference on Artificial Intelligence, **1995**, Montreal, Canada, pp.1137-1143.

28. M. Setnes, R. Babuska, U. Kaymak and H. R. van Nauta Lemke, "Similarity measures in fuzzy rule based simplification", *IEEE Trans. Syst. Man Cybern. Part B Cybern.*, **1998**, *28*, 376-386.

29. C.-L. Huang, M.-C. Chen and C.-J. Wang, "Credit scoring with a data mining approach based on support vector machines", *Expert Syst. Appl.*, **2007**, *33*, 847-856.

30. C. S. Ong, J. J. Huang and G. H. Tzeng, "Building credit scoring models using genetic programming", *Expert Syst. Appl.*, **2005**, *29*, 41-47.

**Appendix**

**Table 9.** Attributes for German credit data set

| No. | Attribute | Type | Value |
|---|---|---|---|
| 1 | Status of existing checking account | Categorical | 0 : ... < 0 DM ; 1 : 0 <= ... < 200 DM<br>2 : ... >= 200 DM /salary assignments for at least 1 year<br>3 : no checking account |
| 2 | Duration in month | Continuous | [4 72] |
| 3 | Credit history | Categorical | 0 : no credits taken/all credits paid back duly; 1 : all credits at this bank paid back duly; 2 : existing credits paid back duly till now; 3 : delay in paying off in the past; 4 : critical account/ other credits existing (not at this bank) |
| 4 | Purpose | Categorical | 0: car (new); 1 : car (used); 2 : furniture/equipment; 3 : radio/television; 4 : domestic appliances; 5 : repairs; 6 : education; 7 : (vacation - does not exist?); 8 : retraining;<br>9 : business; 10 : others |
| 5 | Credit amount | Continuous | [250 18424] |
| 6 | Savings account/bonds | Categorical | 0: ... < 100 DM; 1 : 100 <= ... < 500 DM; 2: 500 <= ... < 1000 DM; 3: .. >= 1000 DM; 4:unknown/ no savings account |
| 7 | Present employment since | Categorical | 0 : unemployed; 1 : ... < 1 years; 2 : 1 <= ... < 4 years ; 3 : 4 <= ... < 7 years; 4 : .. >= 7 years |
| 8 | Installment rate | Continuous | [1 4] |
| 9 | Personal status and sex | Categorical | 0 : male : divorced/separated; 1 : female divorced/ separated /married; 2 : male: single; 3 : male : married/widowed;<br>4 : female : single |
| 10 | Other debtors / guarantors | Categorical | 0 : none; 1 : co-applicant; 3 : guarantor |
| 11 | Present residence since | Continuous | [1 4] |
| 12 | Property | Categorical | 0 : real estate; 1:if not 0 : building society savings agreement/ life insurance; 2: if not 0/1 : car or other, not in attribute 6;<br>3 : unknown / no property |
| 13 | Age in years | Continuous | [19 75] |
| 14 | Other installment plans | Categorical | 0 : bank; 1 : stores; 2 : none |
| 15 | Housing | Categorical | 0 : rent; 1 : own; 2 : for free |
| 16 | Number of existing credits at this bank | Continuous | [1 4] |
| 17 | Job | Categorical | 0 : unemployed/ unskilled - non-resident; 1 : unskilled – resident; 2 : skilled employee / official; 3 : management/ self-employed/ highly qualified employee/ officer |
| 18 | Number of depends | Continuous | [1 2] |
| 19 | Telephone | Categorical | 0 : none; 1 : yes, registered under the customers name |
| 20 | Foreign worker | Categorical | 0 : yes; 1 : no |

**Table 10.** Attributes for Australian credit data set

| No. | Attribute | Type | Value |
|---|---|---|---|
| 1 | A1 | Categorical | 0,1 |
| 2 | A2 | Continuous | [13.75 80.25] |
| 3 | A3 | Continuous | [0 25.125] |
| 4 | A4 | Categorical | 1,2,3 |
| 5 | A5 | Categorical | 1, 2,3,4,5,6 ,7,8,9,10,11, 12, 13,14 |
| 6 | A6 | Categorical | 1, 2,3, 4,5,6,7,8,9 |
| 7 | A7 | Continuous | [0 20] |
| 8 | A8 | Categorical | 1,0 |
| 9 | A9 | Categorical | 1,0 |
| 10 | A10 | Continuous | [0 23] |
| 11 | A11 | Categorical | 1,0 |
| 12 | A12 | Categorical | 1,2,3 |
| 13 | A13 | Continuous | [0 2000] |
| 14 | A14 | Continuous | [1 100001] |