

## Halmstad University

This is an accepted version of a paper published in *Knowledge and Information Systems*. This paper has been peer-reviewed but does not include the final publisher proof-corrections or journal pagination.

Citation for the published paper:

verikas, a., guzaitis, j., gelzinis, a., bacauskiene, m. (2011)

"A general framework for designing a fuzzy rule-based classifier"

*Knowledge and Information Systems*

URL: <http://dx.doi.org/10.1007/s10115-010-0340-x>

Access to the published version may require subscription.

Permanent link to this version:

<http://urn.kb.se/resolve?urn=urn:nbn:se:hh:diva-5825>



<http://diva-portal.org>

# A general framework for designing a fuzzy rule-based classifier

A. Verikas · J. Guzaitis · A. Gelzinis · M. Bacauskiene

Received: July 14 2009 / Revised: December 28 2009 / Accepted: August 22 2010

**Abstract** This paper presents a general framework for designing a fuzzy rule-based classifier. Structure and parameters of the classifier are evolved through a two-stage genetic search. To reduce the search space, the classifier structure is constrained by a tree created using the evolving SOM tree algorithm. Salient input variables are specific for each fuzzy rule and are found during the genetic search process. It is shown through computer simulations of four real world problems that a large number of rules and input variables can be eliminated from the model without deteriorating the classification accuracy. By contrast, the classification accuracy of unseen data is increased due to the elimination.

**Keywords** Classifier, Fuzzy rule, Genetic algorithm, Knowledge extraction, Variable selection, Evolving SOM tree

## 1 Introduction

Support vector machines (SVM) [1], neural networks [2], and relevance vector machines (RVM) [3] are probably the most popular data classification techniques. An SVM and RVM can provide near optimal performance. The advantages of SVM and RVM are the following: the ability to find the global minimum of the objective function, no assumptions made about the data, the complexity of a classifier depends on the number

---

We acknowledge the support from the agency for international science and technology development programmes in Lithuania (COST IC0602).

---

A. Verikas (✉)  
Intelligent Systems Laboratory, Halmstad University,  
Box 823, S-301 18 Halmstad, Sweden  
Tel.: +46 35 167 140  
E-mail: antanas.verikas@hh.se

A. Verikas, J. Guzaitis, A. Gelzinis, M. Bacauskiene  
Department of Electrical & Control Equipment,  
Kaunas University of Technology,  
Studentu 50, 51368, Kaunas, Lithuania  
E-mail: jonas.guzaitis@stud.ktu.lt; adas.gelzinis/marija.bacauskiene@ktu.lt

of support vectors, but not on the dimensionality of the input space. However, in spite of the attempts to explain decisions of such techniques [4, 5], classifiers based on these techniques are not transparent enough and are often considered as “black boxes”. The transparency is very important in some application areas, such as medical decision support or quality control.

By contrast, fuzzy rule-based systems and fuzzy decision trees [6, 7] are known for their transparency and ability of accounting for uncertainty. Fuzzy rule-based classification methods can support rapid incremental learning from new instances without performance degradation on previous training data. ANFIS [8], fuzzy ARTMAP [9], and the fuzzy min-max classifier [10, 11] are examples of the most prominent fuzzy logic-based systems. Fuzzy rule-based systems have been used in a variety of fields such as pattern recognition [12], image segmentation [13], data mining [14], process control [15], resource service selection [16], and system identification [8, 17]. It is well known that designing of fuzzy rule-based systems in high dimensional spaces is rather problematic. However, there are many problems characterized by a small or moderate number of variables. Moreover, quite often high dimensional data vary in a much lower number of dimensions if compared to the dimensionality of the input space. System structure identification and parameter optimization are two main issues to consider when designing a fuzzy rule-based system [8, 18]. Fuzzy partitioning, variable selection, and fuzzy reasoning are the tasks to be solved for identifying the system structure. Parameter optimization usually deals with tuning of parameters of fuzzy membership functions.

Various approaches have been used for dealing with the two main fuzzy rule-based system design issues. The initial system structure, often termed as fuzzy partitioning, is usually identified through K-Means [19], Fuzzy C-Means [20], Learning Vector Quantization (LVQ) [21] or SOM-based clustering [22–28] as well as incremental clustering [29, 30] or by constructing a decision tree [31–33].

Variable selection based on: the output sensitivity to the input change [25, 34], the output sensitivity combined with the correlation between variables [24], Fisher’s inter-class separability measure [35], variable correlation with the output [36] are the most popular variable selection techniques applied. However, quite often, variable selection is not considered at all [26, 27, 23].

It seems that the simple gradient decent [27, 37, 23, 38, 24, 25], error correction [39], and genetic search [26, 35, 31] are the most popular parameter optimization techniques utilized in various studies. The combined optimization of both structure and parameters has also been considered by applying genetic algorithms (GA) [40–44], unsupervised and reinforcement learning [45], or simple heuristics [36]. To reduce the evolution time in the GA-based technique, Chen et al. [46] proposed gathering similar chromosomes into  $k$  clusters and then using a representative chromosome in the evaluation process. In [42, 43] the genetic search process focusses on “hard” data points by assigning a higher weight to such points. Such an approach has also been adopted for learning weights  $z_j^q$  of fuzzy rules [47]. In [48, 49], genetic search-based multi-objective optimization was applied to design a fuzzy rule-based system. The task was to maximize  $f_1(S)$ , minimize  $f_2(S)$ , and minimize  $f_3(S)$ , where  $S$  is a set of fuzzy rules,  $f_1(S)$  stands for correctly classified training samples,  $f_2(S)$  is the number of fuzzy rules in  $S$ , and  $f_3(S)$  is the total number of antecedent conditions in  $S$ . Thus,  $f_3(S)$  can be considered as the total rule length. The optimization starts with all possible rules in the search space defined by the training patterns. In [50], it was proposed to combine several fuzzy rule-based classification systems into a committee through voting or weighted voting.

Generalization ability is an important issue to consider when designing a fuzzy rule-based classifier. The most popular technique applied to improve the generalization ability is rule pruning based on similarity of fuzzy sets [35, 31, 24, 51]. Other approaches utilized are: GA [28], simulated annealing [52], similarity of fuzzy sets combined with GA [26], through forgetting by decaying the grade of certainty of fuzzy rules [39], pruning of rarely used rules [36].

### 1.1 Fuzzy rule-based and nearest neighbour techniques

The fuzzy rule based classification techniques are closely related to nearest neighbour (NN)-based classification approaches. NN-based classification has a sound basis, since there is a considerable body of evidence from the literature that classification and recognition of patterns by humans is best explained as a form of interpolation between similar patterns [53]. NN methods are frequently criticized as requiring much greater use of memory than, for example, neural network algorithms. However, NN learning algorithms can reduce their memory usage by only retaining the full density of exemplars near to classification boundaries and thinning them in other regions [54–57]. As discussed above, various approaches to designing a fuzzy rule-based classifier exploit such techniques for determining an initial system structure through clustering. The location of fuzzy sets, reference patterns in the NN approach, can be further optimized by applying the LVQ techniques.

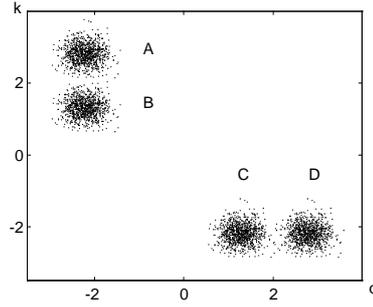
LVQ has been widely used to learn reference patterns for classification based on the NN approach. Each class  $C_j$  is described by several reference patterns  $\mathbf{m}_j^l$  (fuzzy sets in the rule-based approach), which are properly placed within each class region. An unknown  $\mathbf{x}$  is then determined to belong to the class  $k$ , if:

$$k = \arg \min_j [\min_l d(\mathbf{x}, \mathbf{m}_j^l)], \quad l = 1, \dots, N_j, \quad j = 1, \dots, Q \quad (1)$$

where  $Q$  is the number of classes,  $N_j$  is the number of reference patterns representing the class  $j$  and  $d(\mathbf{x}, \mathbf{m}_j^l)$  is the distance between  $\mathbf{x}$  and  $\mathbf{m}_j^l$ .

One more drawback of classical NN and fuzzy rule-based methods is the often exhibited poor generalization performance as compared to neural networks, for example. Neural networks suffer from the "curse of dimensionality" to significantly less extent, since they are able to select useful input features from high-dimensional input vectors. In NN methods, by contrast, the degradation in performance often accompanies the addition of new unimportant features. However, there are many problems where different features are important in different regions of the input space. Fig. 1 provides an example illustrating such a situation.

The four data clusters illustrated in Fig. 1 represent four decision classes. It is obvious that the feature  $q$  is unimportant for discriminating the classes  $A$  and  $B$ , likewise the feature  $k$  is unimportant for discriminating the classes  $C$  and  $D$ . Thus, a subset of features used should be reference pattern or fuzzy rule dependent. However, in most of the known fuzzy rule-based classification algorithms, the feature selection problem is considered independently of the input space region or not considered at all. The objective of this work is to develop a fuzzy modeling framework capable of automatically generating a rule base for classification of numeric data, finding the optimal number of rules and input variables for each rule, and finding the optimal parameter values of fuzzy rules. The remainder of the paper is organized as follows. In



**Fig. 1** Four decision classes in the two-dimensional space.

the next section, the fuzzy model is described. The approach proposed is outlined in Section 3. Section 4 discusses the results of the experimental investigations. Section 5 presents conclusions of the work.

## 2 The fuzzy model

We use the Mamdani model [58], which is the most popular fuzzy model applied in various studies for fuzzy reasoning [26, 25, 23, 59, 24]. Concerning classification, the model is a collection of fuzzy rules  $R_j$  of the following form:

$$R_j : \text{IF } x_1 \text{ is } A_{j1} \text{ AND } \dots \text{ AND } x_n \text{ is } A_{jn} \text{ THEN class } C_q \text{ with } z_j^q \quad (2)$$

where  $A_{ji} (i = 1, \dots, n)$  are fuzzy sets defined over the input variables  $x_i$ ,  $C_q$  is a class label and  $z_j^q$  is a rule weight. Each fuzzy set is represented by a membership function. A triangular or a Gaussian function of the form

$$\mu_{ji} = \exp\left(-\frac{(x_i - c_{ji})^2}{\sigma_{ji}^2}\right) \quad (3)$$

where  $c_{ji}$  and  $\sigma_{ji}$  are the center and the width of the Gaussian function, respectively, are common choices. We use Gaussian membership functions in this study. There are various ways to determine the rule weights  $z_j^q$ . In this work, the weight  $z_j^q$  is given by:

$$z_j^q = \max\left(\frac{\sum_{\mathbf{x}_p \in C_q} \mu_{\mathbf{A}_j}(\mathbf{x}_p) - \sum_{\mathbf{x}_p \notin C_q} \mu_{\mathbf{A}_j}(\mathbf{x}_p)}{\sum_{p=1}^N \mu_{\mathbf{A}_j}(\mathbf{x}_p)}, 0\right) \quad (4)$$

where  $N$  is the number of training patterns and the matching degree of the input pattern  $\mathbf{x}_p$  with the antecedent part  $\mathbf{A}_j = (A_{j1}, \dots, A_{jn})$  is calculated using a  $T$ -norm

$$\mu_{\mathbf{A}_j}(\mathbf{x}_p) = T(\mu_{A_{j1}}(x_{p1}), \dots, \mu_{A_{jn}}(x_{pn})) \quad (5)$$

We use the *min*  $T$ -norm operator in this work. Weights  $z_j^q$  of this type were studied in [48]. In [60], weights based on ROC analysis are advocated.

A winning rule is used to make a decision. Thus, given a rule base  $S$  consisting of  $L$  rules, an input pattern  $\mathbf{x}_p$  is assigned to the class  $q$  if

$$q = \arg \max_k \{T[\mu_{\mathbf{A}_j}(\mathbf{x}_p), z_j^k], j = 1, \dots, L\} \quad (6)$$

where  $T$  is the *product*  $T$ -norm operator, in this work.

Having defined the membership functions, we formulate the fuzzy modeling problem in the following way. Given  $N$  pairs of input-output patterns  $(\mathbf{x}, y)$ , create a minimal number of fuzzy rules  $r$  with the optimal number of features  $n_i$  for each rule and find the optimal values of parameters  $(\mathbf{c}, \boldsymbol{\sigma}, \mathbf{z})$  of the fuzzy model  $F(\mathbf{c}, \boldsymbol{\sigma}, \mathbf{z}, r, n_i, i = 1, \dots, r)$ .

### 3 The approach

The procedure to construct the fuzzy rule-based classifier consists of the following steps.

1. Divide the data set into learning and test subsets.
2. Cluster the learning set data by applying the evolving SOM tree.
3. Based on the evolved tree, generate a population of sub-trees. Each sub-tree defines the initial structure of one fuzzy rule-based classifier. The generation is accomplished by randomly cutting branches of the tree grown in Step 2. The cutting occurs approximately between 25 and 75% of the tree depth.
4. Represent each node in the sub-tree population by a set of fuzzy sets with the Gaussian membership functions.
5. Take one sub-tree (classifier of a given structure) from the sub-tree population and encode the structure, features (used/not used), and parameters of the membership functions of the classifier into a chromosome. When encoding, enable feature selection independently for each fuzzy rule.
6. Generate a population of chromosomes encoding individual classifiers of the given structure. The individual classifiers differ in features and values of the parameters.
7. Apply the modified LVQ-3 algorithm to the individuals of the population.
8. Evaluate the fitness of the individuals.
9. Apply genetic operations (to features and parameters) and generate a new population.
10. Repeat Steps 7–9 until convergence.
11. Take the best individual of the given structure.
12. Repeat Steps 5–11 for the whole population of sub-trees.
13. Apply genetic operations (to structure of sub-trees) and generate a new population of sub-trees.
14. Repeat Steps 5–13 for a given number of generations.

Next, we briefly describe the main topics of the technique.

### 3.1 The algorithm

1. create evolving SOM tree
2. generate and encode a population of sub-trees  $tp$
3. **for** (given number of generations): *structure evolution*
4.     **for** (each chromosome  $tc$  in  $tp$ ): *parameter evolution*
5.         generate a population  $pp$  encoding individual classifiers of the  $tc$  structure
6.         **repeat**
7.             **for** (each chromosome  $pc$  in  $pp$ )
8.                 apply the modified LVQ-3 algorithm to  $pc$
9.                 evaluate the fitness of  $pc$
10.             **end for**
11.             apply genetic operations to  $pp$
12.             **until** (convergence)
13.             take the best individual of  $pp$  for  $tc$  in  $tp$
14.         **end for**
15.         apply genetic operations to  $tp$
16.     **end for**

### 3.2 The evolving SOM tree

Like SOM, the evolving SOM tree [61] exhibits the self-organization property. The evolving tree structure enables the SOM tree to efficiently handle large scale problems. Moreover, there is no need of choosing the map size beforehand. Like in ordinary SOM, each node of the SOM tree has a weight vector  $\mathbf{w}_i$ . When training the tree, for each training vector  $\mathbf{x}$ , the best matching unit (BMU) is found by a greedy tree search. BMU is always a leaf node. Weight vectors of the BMU and its neighbours are then updated using the SOM adaptation rule:

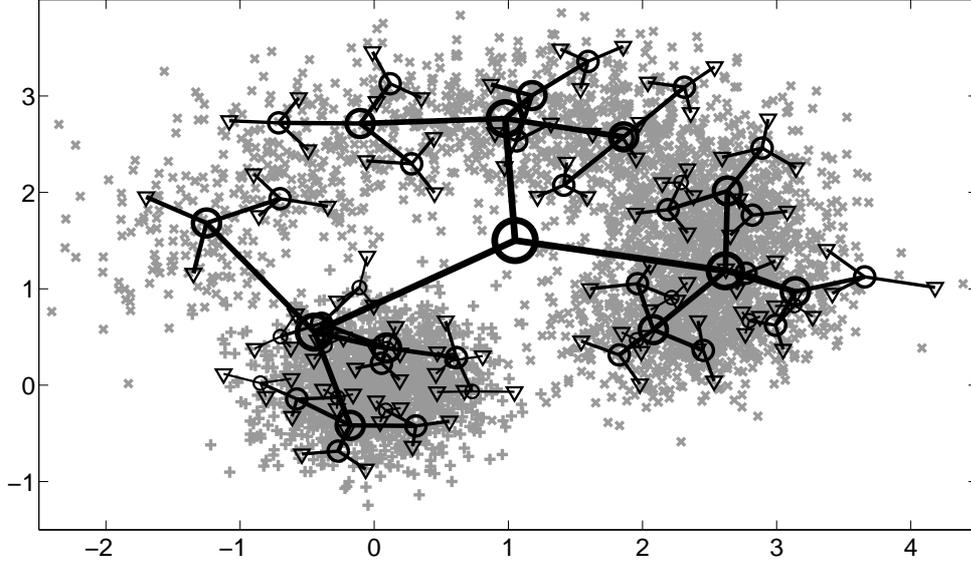
$$\mathbf{w}_i(t+1) = \mathbf{w}_i(t) - h_{ci}(t)[\mathbf{x}(t) - \mathbf{w}_i(t)] \quad (7)$$

where  $h_{ci}(t)$  is the neighbourhood function. We used the Gaussian neighbourhood function

$$h_{ci}(t) = \beta(t) \exp\left(-\frac{\|\mathbf{r}_c - \mathbf{r}_i\|^2}{2s^2(t)}\right) \quad (8)$$

where  $s(t)$  is the width of the Gaussian function,  $\mathbf{r}_c$  and  $\mathbf{r}_i$  denote location of nodes  $c$  and  $i$ , and  $\beta(t)$  is the learning rate. The meaning of  $s(t)$  and  $\beta(t)$  is the same as in SOM [62], while the meaning of the norm  $\|\mathbf{r}_c - \mathbf{r}_i\|$  is quite different. The basic idea of calculating the distance is to count how many "hops" are needed to get from the BMU to the considered node along the shortest path [61]. The distance  $\|\mathbf{r}_c - \mathbf{r}_i\|$  is then given by the number of hops minus one.

Fig 2 presents an example of the evolving SOM tree generated to represent the two-dimensional data, where circles denote the center points of the SOM nodes. The deeper is a node in the tree, the smaller circle and the thinner line are used in the visualization. The same applies for the links. The leaf nodes are denoted by triangles. Each parent node is split into three child nodes, in the example. However, other number of child nodes can be used. Each node is characterized by a center point and width values computed from the data mapped onto the node. These parameters are used as initial values for the parameters of the Gaussian membership functions.



**Fig. 2** An example of the evolving SOM tree generated to represent the two-dimensional data.

### 3.3 The modified LVQ-3 algorithm

Assume that  $d_i$  and  $d_j$  are the Euclidean distances from the pattern  $\mathbf{x}$  to the reference patterns  $\mathbf{m}_i$  and  $\mathbf{m}_j$ , respectively. Then  $\mathbf{x}$  is defined to fall into a window of the relative width  $\lambda$ , if

$$\min\left(\frac{d_i}{d_j}, \frac{d_j}{d_i}\right) > \frac{1-\lambda}{1+\lambda} \quad (9)$$

For all  $\mathbf{x}$  falling into the window adapt:

$$\mathbf{m}_i(t+1) = \mathbf{m}_i(t) - \alpha(t)[\mathbf{x}(t) - \mathbf{m}_i(t)] \quad (10)$$

$$\mathbf{m}_j(t+1) = \mathbf{m}_j(t) + \alpha(t)[\mathbf{x}(t) - \mathbf{m}_j(t)] \quad (11)$$

where  $\alpha(t)$  decreases with time and  $0 < \alpha(t) < 1$ ,  $\mathbf{m}_i$  and  $\mathbf{m}_j$  are two closest reference patterns to  $\mathbf{x}$ , whereby  $\mathbf{x}$  belongs to the same class as  $\mathbf{m}_j$ , but not as  $\mathbf{m}_i$ . If  $\mathbf{x}$ ,  $\mathbf{m}_i$  and  $\mathbf{m}_j$  belong to the same class:

$$\mathbf{m}_k(t+1) = \mathbf{m}_k(t) + \gamma\alpha(t)[\mathbf{x}(t) - \mathbf{m}_k(t)] \quad (12)$$

for  $k \in \{i, j\}$ . If  $\mathbf{x}$  belongs to a different class than  $\mathbf{m}_i$  and  $\mathbf{m}_j$ :

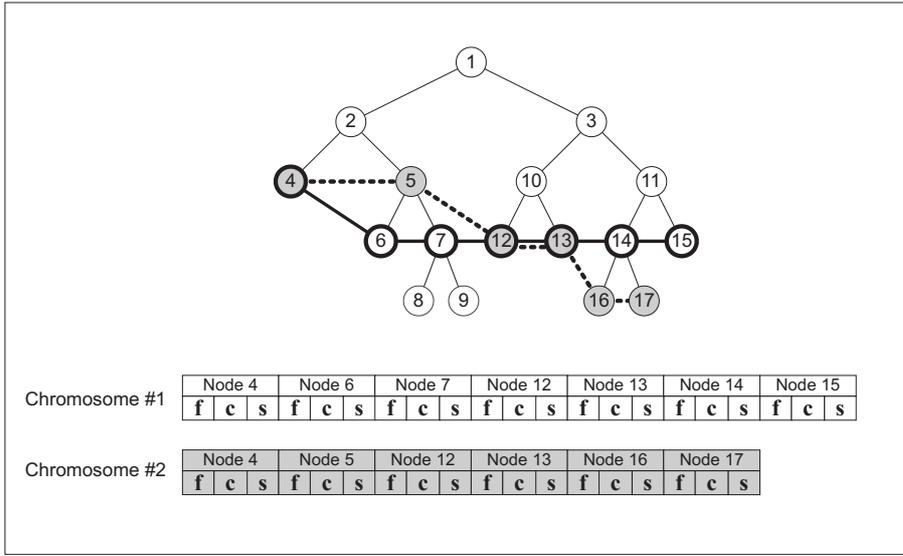
$$\mathbf{m}_k(t+1) = \mathbf{m}_k(t) - \gamma\alpha(t)[\mathbf{x}(t) - \mathbf{m}_k(t)] \quad (13)$$

for  $k \in \{i, j\}$ . The optimal value of the parameter  $\gamma$  depends on the size of the window, the value is smaller for narrower windows [21]. Values between 0.1 and 0.5 are suggested for  $\gamma$  [21].

The algorithm performs fine tuning of the centers of membership functions and helps reducing the time of genetic search. The last adaptation step is not used in the original version of the Lvq-3 algorithm. We have found that the use of the step quite noticeably improved the accuracy of the algorithm.

### 3.4 Encoding

Structure, features, and parameters of the membership functions are to be encoded. The structure is determined by a sub-tree and is encoded as a connected graph. Fig. 3 presents an example of the evolving SOM tree along with two chromosomes encoding two hypothetical sub-trees. The hypothetical leaf nodes of the two sub-trees are shown connected by the bold solid and the dashed line, respectively, in Fig. 3. There are as many *sections* in the chromosome, as there are leaf nodes in the corresponding sub-tree (the number of fuzzy rules in the classifier).

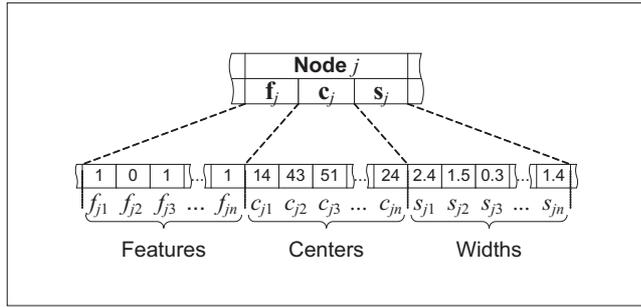


**Fig. 3** An example of the evolving SOM tree along with two chromosomes encoding two hypothetical sub-trees, the leaf nodes of which are connected by the bold solid and dashed lines.

Each chromosome section consists of three *sub-sections*: **f**, **c**, and **s**, encoding features ("feature mask"), centers and widths of the membership functions  $\mu$ , respectively. Fig. 4 illustrates sub-sections of the  $j$ th node. There are  $n$  bits in the  $\mathbf{f}_j$  sub-section, where  $n$  is the dimensionality of the input space. A bit in the  $\mathbf{f}_j$  sub-section set to 0/1 means that the corresponding feature is *used/not used* in the rule encoded in the corresponding section of the chromosome. Centers and widths of the membership functions  $\mu$  are stored as real/integer numbers in slots of the subsections  $\mathbf{c}_j$  and  $\mathbf{s}_j$ . For example, integer numbers can be used to encode centers of the membership functions in the applications related to image analysis. There are  $n$  slots in each,  $\mathbf{c}_j$  and  $\mathbf{s}_j$ , sub-section.

### 3.5 Genetic operations

Crossover and mutation are the genetic operations applied in both loops of genetic evolution: the loop concerning structure evolution and the loop concerning features

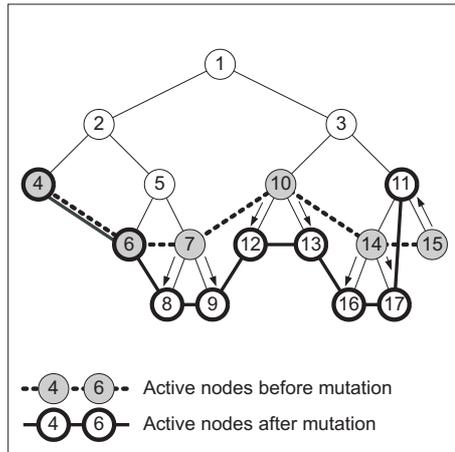


**Fig. 4** Three sub-sections ( $f_j$ ,  $c_j$ , and  $s_j$ ) of the  $j$ th node encoding features, centers and width of the membership functions, respectively.

and parameters of the membership functions. The crossover and mutation operations are executed with the probability of crossover  $p_c$  and the mutation probability  $p_m$ , respectively.

### 3.5.1 Mutation of chromosomes encoding structure

Mutation in structure evolution amounts to taking one step up or down (the direction is selected randomly) along a randomly selected branch of the tree. To select the direction, an integer  $\vartheta$  is selected randomly from the set  $\{-1, 1\}$  for each node. The node undergoing the mutation is replaced with its parent node, if  $\vartheta < 0$  or is split into children, if  $\vartheta > 0$ . Fig. 5 illustrates the mutation operation in structure evolution.



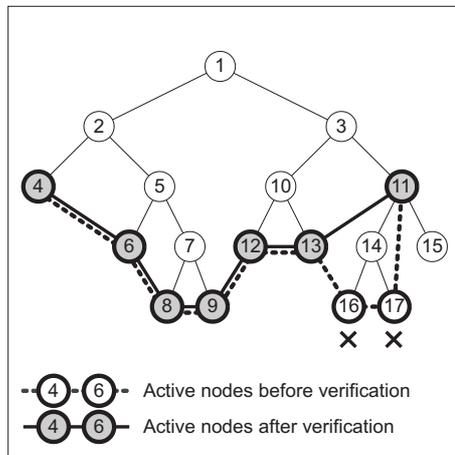
**Fig. 5** A part of the evolving SOM tree illustrating the mutation operation in structure evolution.

Nodes being active before the mutation operation are shown in grey and connected by the dashed line, in Fig. 5. A node is said to be active if it is used to define the classifier structure at a current moment. Nodes being active after the mutation operation are shown connected by the bold line, in Fig. 5. Arrows in Fig. 5 indicate the direction of

mutation. As can be seen, nodes 7, 10, and 14 mutate towards their children, node 15 mutates towards its parent, while nodes 4 and 6 do not mutate.

### 3.5.2 Verifying the structure consistency

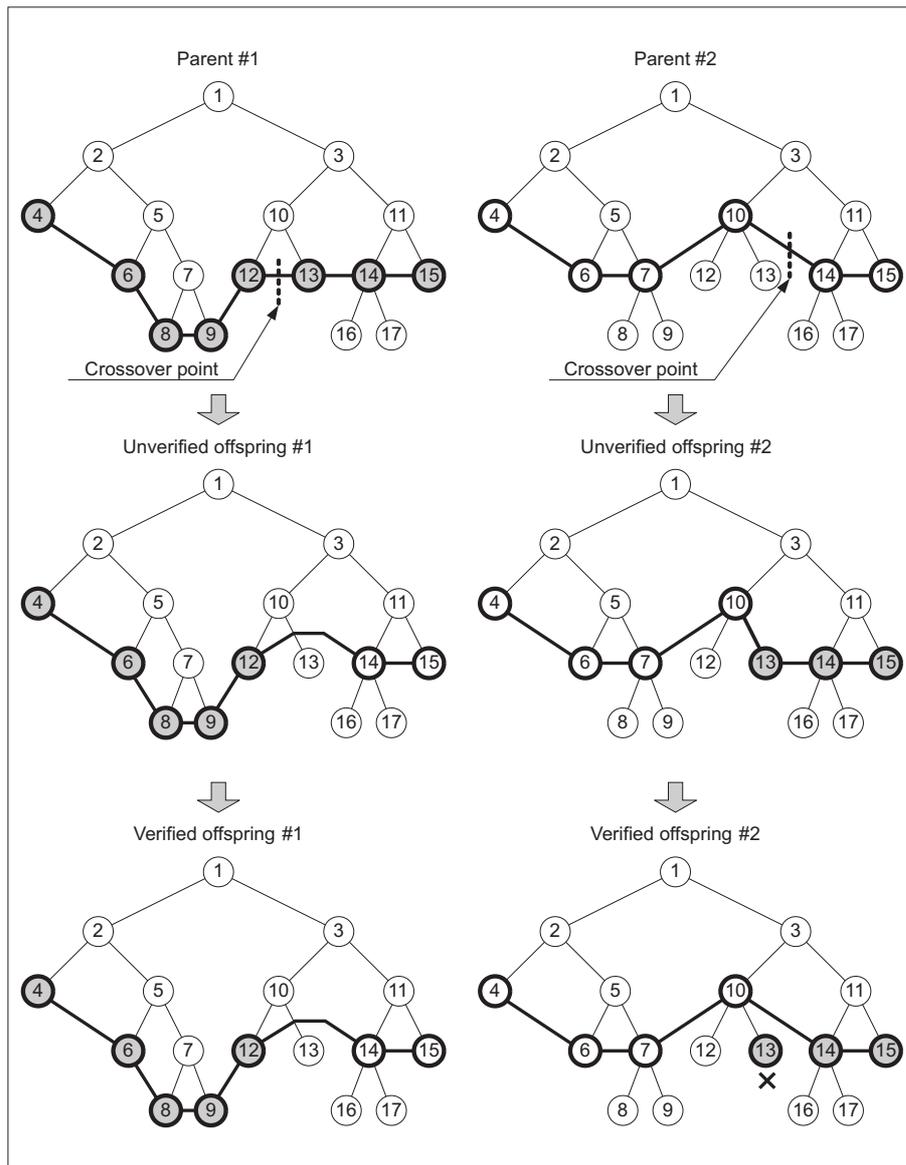
After the genetic operations, the sub-trees are verified for consistency. The redundant active nodes are deactivated during the verification. Fig. 6 illustrates the verification operation. Nodes being active before the verification operation, are shown in grey and connected by the dashed line, in Fig. 6. Nodes being active after the verification operation are shown connected by the bold line. As can be seen from Fig. 6, nodes 4, 6, 8, 9, 12, and 13 are left unaffected, while nodes 16 and 17 are deactivated, since node 11, predecessor of nodes 16 and 17, was activated during the mutation operation. The deactivated nodes are labelled by  $\times$ , in Fig. 6.



**Fig. 6** A part of the evolving SOM tree illustrating the verification operation in structure evolution.

### 3.5.3 Crossover of chromosomes encoding structure

When performing crossover for structure evolution, parts of two sub-trees are exchanged. The crossover point is selected randomly. Fig. 7 illustrates two sub-trees (*Parent #1* and *Parent #2*) before the crossover operation with the crossover points indicated by the dashed lines. The active nodes of the trees encoded in the corresponding chromosomes are shown connected by the bold lines. The resulting sub-trees after the crossover operation are labelled as *Unverified offspring #1* and *Unverified offspring #2*, in Fig. 7. The verification procedure is then applied to the unverified offsprings. The verification results into the verified offsprings with one deactivated node, namely node 13 of the offspring #2.

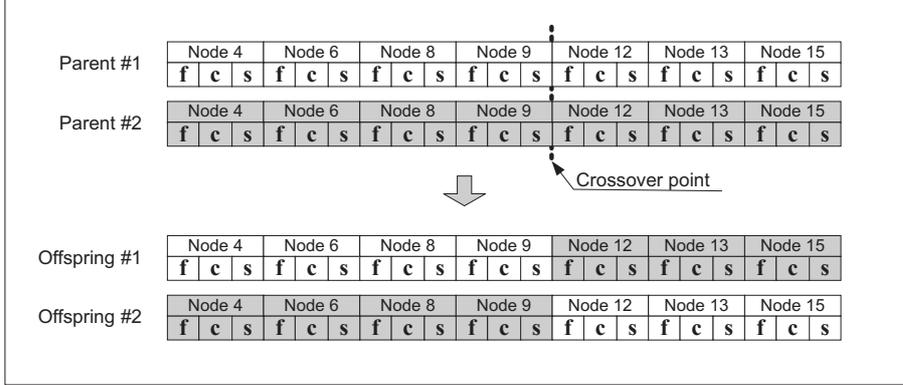


**Fig. 7** A part of the evolving SOM tree illustrating the crossover operation in structure evolution.

### 3.5.4 Crossover and mutation of chromosomes encoding features and parameters

Crossover of two chromosomes encoding features and parameters of the membership functions are performed at a randomly selected point by exchanging parts of the chromosomes, as it is illustrated in Fig. 8. Two chromosomes selected for the crossover operation are to encode the same number of nodes.

The mutation operation is accomplished by reversing the value of a bit in the “feature mask” (**f** sub-section) and by adding a random value from a given, symmetric interval to parameters stored in the slot (gene) selected for mutation in the **c** and **s** sub-sections. Selection of genes for mutation is performed independently in the three sub-sections.



**Fig. 8** Parameters of two parent chromosomes are exchanged at the crossover point.

### 3.6 Fitness function

The fitness value of the  $i$ th chromosome  $f_i$  is given by

$$f_i = \chi_i - \eta \frac{\sum_{j=1}^{r_i} n_j}{r_0 \times n} \quad (14)$$

where  $\chi_i$  is the classification accuracy obtained from the classifier encoded in the  $i$ th chromosome,  $\eta$  is a parameter determining the degree to which classifiers with a large number of features are penalized,  $n$  is the total number of available features,  $n_j$  is the number of features used by the  $j$ th rule,  $r_0$  is the number of rules in the initial tree, and  $r_i$  stands for the number of rules used by the  $i$ th classifier. A value of  $\eta = 0.1$  worked well in all the tests.

Chromosomes are selected for genetic operations with some probability. The selection probability of the  $i$ th chromosome  $p_i$  is given by

$$p_i = \frac{f_i}{\sum_{j=1}^M f_j} \quad (15)$$

where  $M$  is the population size. The roulette selection principle was applied.

## 4 Experimental Investigations

### 4.1 Data Used

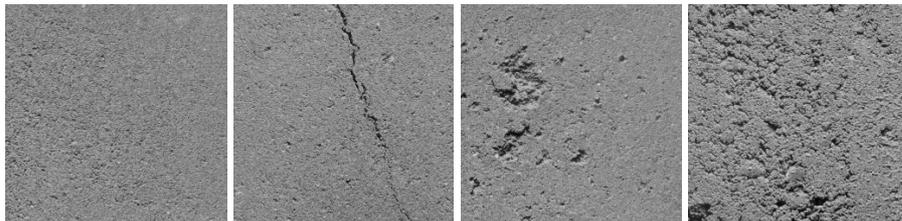
Four data sets have been used in the tests.

**US congressional voting records problem.** The United States Congressional voting records data set consists of the voting records of 435 congressman on 16 major issues in the 98th Congress. The votes are categorized into one of the three types of votes: (1) (*Yea*), (2) (*Nay*), and (3) (*Unknown*). The task is to predict the correct political party affiliation of each congressman. The 98th Congress consisted of 267 Democrats and 168 Republicans.

**Wisconsin diagnostic breast cancer problem.** There are 30 real-valued features. The features are computed from a digitized image of a fine needle aspirate of a breast mass and describe characteristics of the cell nuclei present in the image. There are 569 instances, 357 benign and 212 malignant.

**The diabetes diagnosis problem.** The *Pima Indians Diabetes* data set contains 768 samples taken from patients who may show signs of diabetes. Each sample is described by eight features: (1) *Number of times pregnant*, (2) *Plasma glucose concentration*, (3) *Diastolic blood pressure*, (4) *Triceps skin fold thickness*, (5) *Two-hour serum insulin*, (6) *Body mass index*, (7) *Diabetes pedigree function*, and (8) *Age*. There are 500 samples from patients who do not have diabetes and 268 samples from patients who are known to have diabetes. These three data sets are available at: <http://archive.ics.uci.edu/ml/>.

**Pavement tiles surface inspection problem.** A pavement tile surface is to be assigned into a *quality* or *defective* class. Features for the classification are extracted from a camera image. Five features characterizing the image texture and the grey level distribution [63] have been used to design a classifier. Fig. 9 presents four examples of pavement tile surfaces used in the study. In total, 200 quality and 200 defective surfaces were available.



**Fig. 9** Examples of pavement tile surfaces: a quality surface on the left and three defective surfaces.

## 4.2 Experimental Setup

We randomly assign the available data points into the learning  $D_L$  and test  $D_T$  data sets. The data are normalized to have zero mean and unit variance. We run an experiment 30 times with different random splits into the sets  $D_L$  and  $D_T$ . The results obtained are averaged over the 30 runs.

### 4.3 Optimization Parameters

The optimal size of the Lvq-3 window depends on the number of training samples. If a large number of samples is available, a narrow window would guarantee the most accurate location of the decision boundary. For good statistical accuracy, however, the number of samples falling into the widow must be sufficient [21]. The optimal value of  $\gamma$  depends on the size of the window, being smaller for narrower windows [21,64]. After some experiments the following values of the Lvq-3 parameters have been used:  $\lambda = 0.05$ ,  $\alpha = 0.02$ , and  $\gamma = 0.4$ .

There are two loops of genetic evolution: the outer loop concerning structure evolution and the inner loop concerning features and parameters of the membership functions. The genetic search lasted for 100 generations (for both loops) with the following parameters: the population size was set to 50 and the number of offsprings produced for creating the next population was equal to 50. The number of generations was determined experimentally by monitoring changes of the fitness function value. The number chosen was such that no fitness function value increase was observed in the last 10 generations. The values of crossover and mutation probabilities were found experimentally. The following values worked well in the tests:  $p_c = 0.95$  and  $p_m = 0.01$ . The appropriate value of the parameter  $\eta$  was found to be  $\eta = 0.05$ .

### 4.4 Results

In the first set of experiments, the feature selection has not been applied and the classification accuracy obtained from the fuzzy rule-based classifier was compared with the accuracy achieved by other techniques. The multi layer perceptron (MLP), k-NN, and LVQ-3 classifiers have been used for the comparison. The appropriate number of hidden nodes in the MLP and the  $k$  value of the k-NN classifier were found experimentally. The leaf nodes created by the evolving SOM tree were used as initial reference patterns for the LVQ-3 classifier. Table 1 presents the average test data set classification accuracy (%) obtained from the different classifiers using all available features. As can be seen from Table 1, the proposed fuzzy rule-based classifier provided the highest classification accuracy for all the problems studied.

**Table 1** The average test data set classification accuracy (%) obtained from different classifiers using all available features.

Data Set\Classifier	k-NN	MLP	LVQ-3	Proposed
Voting	91.24	93.78	77.41	94.68
Breast cancer	88.73	97.18	75.87	98.54
Diabetes	71.19	71.49	64.22	71.92
Surface inspection	77.30	81.63	78.13	84.13

In the next set of experiments, feature selection was activated and features, specific for each rule were found through the genetic search. Table 2 presents the average test data set classification accuracy obtained from the approach proposed using the selected features. The classification accuracy obtained using all the available features

is also presented for the sake of comparison. The obtained improvement in classification accuracy should be obvious from Table 2. Assuming that the classification errors are log-normally distributed and applying the  $t$ -test it was found that the difference between the classification accuracy obtained using the selected and all features is significant with 95% confidence, except for the Breast cancer data. In Table 2, we also provide the average test data set classification accuracy obtained by other authors for the same public data sets in recent studies. In [65–67], evolutionary techniques have also been used to design the classifiers. As can be seen from Table 2, the proposed technique outperformed the other approaches on all the three data sets.

**Table 2** The average test data set classification accuracy (%) obtained from the approach proposed using all and selected features, along with the average accuracy obtained by other authors in recent studies.

Data Set\Features	All features	Selected features	Comparison
Voting	94.68	98.68	96.98 [65], 95.31 [66]
Breast cancer	98.54	99.02	95.55 [67], 94.38 [47]
Diabetes	71.92	75.92	75.19 [67], 75.00 [60]
Surface inspection	84.13	99.63	—

Table 3 presents information on the number of rules and features used to classify the data. In the parentheses given are the number of initial rules and the number of available features. Ranges in the “Features” column indicate the minimum and the maximum number of features used by different rules. As can be seen from Table 3, the number of features used by different rules varies significantly. Observe that even if the number of features used by two different rules is the same, the features used are often different. Thus, features used are rule specific, indeed.

**Table 3** The number of rules and features used to classify data from the different data sets.

Data Set	# Rules	# Features
Voting	15 (25)	05–09 (16)
Breast cancer	10 (14)	12–20 (30)
Diabetes	10 (16)	02–08 ( 8)
Surface inspection	09 (13)	03–04 ( 5)

Below given is an example of a fuzzy rule (one out of ten) generated by the proposed technique for classification of the Diabetes set data.

$R_1$  : IF  $x_2$  is ABOUT 119.0 AND  $x_4$  is ABOUT 38.1 THEN class *Healthy* with 0.84

where “ABOUT” is a fuzzy variable, 119.0 and 38.1 are the center points of the Gaussian membership functions,  $x_2$  stands for “plasma glucose concentration”, and  $x_4$  means “triceps skin fold thickness”. Knowing meaning of the variables  $x_1, \dots, x_8$ , the rules are easy to interpret for a medical doctor.

One may wonder what classification accuracy would be achieved, if an ordinary GA-based feature selection procedure—not specific for each fuzzy rule—were used.

Table 4 presents the results obtained from such an experiment. Comparing the results presented in Table 2, Table 3, and Table 4 we can see that the fuzzy rule specific feature selection procedure results into the decreased average number of rules and features and the increased classification accuracy. While the difference between the average number of rules and features utilized by the two techniques is rather small, the difference between the classification accuracies is statistically significant with 95% confidence for all the data sets.

**Table 4** The average number of rules and features used to classify data from the different data sets and the average test data set classification accuracy (%) obtained using an ordinary GA-based feature selection procedure.

Data Set	# Rules	# Features	Accuracy
Voting	16 (25)	08 (16)	95.12
Breast cancer	10 (14)	17 (30)	96.91
Diabetes	11 (16)	06 ( 8)	73.23
Surface inspection	10 (13)	04 ( 5)	98.11

Next, the influence of crossover and mutation probabilities,  $p_c$  and  $p_m$ , on classification accuracy was studied. The same  $p_c$  and  $p_m$  values were used for both structure and parameter evolution. To reduce the computation time and to decouple the influence of  $p_c$  and/or  $p_m$ , and feature selection on the classification accuracy, the studies were performed without employing feature selection. A similar performance was observed for  $p_m$  values ranging from 0.005 to 0.05. A value of  $p_m = 0.01$  was selected. When studying the influence of  $p_c$ ,  $p_m$  was set to  $p_m = 0$ . Table 5 presents the average test set classification accuracy obtained for different  $p_c$  values. The interval of  $p_c = [0.8, 1.0]$  was studied additionally and was found that the highest and similar classification accuracy is obtained for  $p_c = 0.9$  and  $p_c = 0.95$ . Thus,  $p_c$  values close to unity are recommended.

**Table 5** The average test set classification accuracy (%) obtained for different  $p_c$  values.

Data Set	$p_c = 0.25$	$p_c = 0.50$	$p_c = 0.75$	$p_c = 1.0$
Voting	92.89	93.01	93.21	94.68
Breast cancer	94.12	93.32	94.78	95.02
Diabetes	67.21	68.45	69.53	71.92
Surface inspection	82.01	82.10	83.02	84.13

One and the same loop can be used for both structure and parameter evolution of the classifier, presumably at the expense of computation time. An experiment was performed to compare these two implementations using the *Diabetes* data set. A very similar test data set classification accuracy was achieved in both implementations. However, the evolution based on the two-loop implementation was about five times faster.

## 5 Conclusions

Proposed is a general framework for designing a fuzzy rule-based classification system. The developed two-stage GA partitions the search space and enables evolving both structure and parameters of the classifier. Salient input variables, specific for each fuzzy rule, are also found during the search process.

Computer simulations of four real world problems have shown that the performance obtained from the classifier is comparable or even higher than the best performance obtained by other authors when using “black box” as well as fuzzy rule-based models. The proposed variable selection tool allowed to significantly increase the classification accuracy if compared to the case of using all the available input variables. It was shown through computer simulations that a large number of rules and input variables can be eliminated from the model without deteriorating the classification accuracy. Moreover, the classification accuracy of the test set data increased due to the reduction.

## References

1. Vapnik, V.N.: *Statistical Learning Theory*. Wiley, New York (1998)
2. Bishop, C.M.: *Pattern Recognition and Machine Learning*. Springer, Singapore (2006)
3. Tipping, M.E.: Sparse bayesian learning and the relevance vector machine. *Journal of Machine Learning Research* **1** (2001) 211–244
4. Hamel, L.: Visualization of support vector machines with unsupervised learning. In: *Proceedings of IEEE Symposium Computational Intelligence and Bioinformatics and Computational Biology, 2006 CIBCB'06, IEEE* (2006) 1–8
5. Strumbelj, E., Bosnic, Z., Kononenko, I., Zakotnik, B., Kuhar, C.G.: Explanation and reliability of prediction models: the case of breast cancer recurrence. *Knowledge and Information Systems* **24**(2) (2010) 305–324
6. Janikow, C.Z.: Fuzzy decision trees: issues and methods. *IEEE Trans Systems Man, Cybernetics–Part B: Cybernetics* **28**(1) (1998) 1–14
7. Olaru, C., Wehenkel, L.: A complete fuzzy decision tree technique. *Fuzzy Sets and Systems* **138** (2003) 221–254
8. Jang, J.R.: ANFIS: adaptive-network-based fuzzy inference system. *IEEE Trans Systems Man Cybernetics* **23** (1993) 665–685
9. Carpenter, G.A., Grossberg, S., Markuzon, N., Reynolds, J.H., Rosen, D.B.: Fuzzy ARTMAP: A neural network architecture for incremental supervised learning of analog multidimensional maps. *IEEE Trans Neural Networks* **3**(5) (1992) 698–713
10. Simpson, K.P.: Fuzzy min-max neural networks - Part 1: Classification. *IEEE Trans Neural Networks* **3**(5) (1992) 776–786
11. Abe, S., Lan, M.S.: A method for fuzzy rules extraction directly from numerical data and its application to pattern classification. *IEEE Trans Fuzzy Systems* **3**(1) (1995) 18–28
12. Pedrycz, W.: Fuzzy sets in pattern recognition: Methodology and methods. *Pattern Recognition* **23**(1-2) (1990) 121–146
13. Jaffar, M.A., Hussain, A., Mirza, A.M.: Fuzzy entropy based optimization of clusters for the segmentation of lungs in CT scanned images. *Knowledge and Information Systems* **24**(1) (2010) 91–111
14. Denton, A.M., Wu, J.: Data mining of vector-item patterns using neighborhood histograms. *Knowledge and Information Systems* **21** (2009) 173–199
15. Takagi, T., Sugeno, M.: Fuzzy identification of systems and its applications to modeling and control. *IEEE Trans Systems, Man Cybernetics* **15**(1) (1985) 116–132
16. Tao, F., Zhao, D., Zhang, L.: Resource service optimal-selection based on intuitionistic fuzzy set and non-functionality QoS in manufacturing grid system. *Knowledge and Information Systems* (2009) DOI 10.1007/s10115-009-0263-6
17. Wang, L.X.: Modeling and control of hierarchical systems with fuzzy systems. *Automatica* **33**(6) (1997) 1041–1053
18. Nauck, D., Klawon, F., Kruse, R.: *Foundations of Neuro-Fuzzy Systems*. Wiley, Chichester (1997)

19. Duda, R.O., Hart, P.E., Stork, D.G.: Pattern Classification. 2 edn. John Wiley & Sons, New York (2001)
20. Bezdek, J.C.: Pattern Recognition with Fuzzy Objective Function Algorithms. Plenum Press, New York (1981)
21. Kohonen, T.: The self-organizing map. *Proceedings of the IEEE* **78**(9) (1990) 1461–1480
22. Yoshinari, Y., Pedrycz, W., Hirota, K.: Construction of fuzzy models through clustering techniques. *Fuzzy Sets and Systems* **54**(2) (1993) 157–165
23. Wang, Y., Rong, G.: A self-organizing neural-network based fuzzy system. *Fuzzy Sets and Systems* **103** (1999) 1–11
24. Chen, M.Y., Linkens, D.A.: A systematic neuro-fuzzy modeling framework with application to material property prediction. *IEEE Trans Systems Man, Cybernetics, Part B* **31**(5) (2001) 781–790
25. Castellano, G., Castiello, C., Fanelli, A.M., Mencar, C.: Knowledge discovery by a neuro-fuzzy medeling framework. *Fuzzy Sets and Systems* **149** (2005) 187–207
26. Zhou, E., Khotanzad, A.: Fuzzy classifier design using genetic algorithms. *Pattern Recognition* **40**(12) (2007) 3401–3414
27. Nishina, T., Hagiwara, M.: Fuzzy inference neural network. *Neurocomputing* **14** (1997) 223–239
28. Chang, P.C., Liao, T.W.: Combining som and fuzzy rule base for flow time prediction in semiconductor manufacturing factory. *Applied Soft Computing* **6** (2006) 198–206
29. Kasabov, N.: Evolving fuzzy neural networks for supervised/unsupervised on-line, knowledge-based learning. *IEEE Trans Systems, Man and Cybernetics* **31**(6) (2001) 902–918
30. Minku, F.L., Ludermir, T.B.: Clustering and co-evolution to construct network ensembles: An experimental study. *Neural Networks* **21** (2008) 1363–1379
31. Abonyi, J., Roubos, J.A., Szeifert, F.: Data-driven generation of compact, accurate, and linguistically sound fuzzy classifiers based on a decision-tree initialization. *International Journal of Approximate Reasoning* **32**(1) (2003) 1–21
32. Kbir, M.M., Benkirane, H., Maalmi, K., Benslimane, R.: Hierarchical fuzzy partition for pattern classification with fuzzy if-then rules. *Pattern Recognition Letters* **21** (2000) 503–509
33. Pulkkinen, P., Koivisto, H.: Identification of interpretable and accurate fuzzy classifiers and function estimators with hybrid methods. *Applied Sof Computing* **7** (2007) 520–433
34. Verikas, A., Bacauskiene, M.: Feature selection with neural networks. *Pattern Recognition Letters* **23**(11) (2002) 1323–1335
35. Roubos, J.A., Setnes, M., Abonyi, J.: Learning fuzzy classification rules from labeled data. *Information Sciences* **150** (2003) 77–93
36. Nauck, D., Kruse, R.: Obtaining interpretable fuzzy classification rules from medical data. *Artificial Intelligence in Medicine* **16**(2) (1999) 149–169
37. Lin, C.T., Lu, Y.C.: A neural fuzzy system with linguistic teachnig signals. *IEEE Trans Fuzzy Systems* **3**(2) (1995) 169–189
38. Tung, W.L., Quek, C., Cheng, P.: GenSo: a novel neural-fuzzy based early warning system for predicting bank failures. *Neural Networks* **17** (2004) 567–587
39. Nozaki, K., Ishibuchi, H., Tanaka, H.: Adaptive fuzzy rule-based classification systems. *IEEE Trans Fuzzy Systems* **4**(3) (1996) 238–250
40. Wang, C.H., Hong, T.P., Tseng, S.S.: Integrating fuzzy knowledge by genetic algorithms. *IEEE Transactions on Evolutionary Computation* **2**(4) (1998) 138–149
41. Tsang, C.H., Kwong, S., Wang, H.: Genetic-fuzzy rule mining approach and evaluation of feature selection techniques for anomaly intrusion detection. *Pattern Recognition* **40**(9) (2007) 2373–2391
42. Ozyer, T., Alhadj, R., Barker, K.: Intrusion detection by integrating boosting genetic fuzzy classifier and data mining criteria for rule pre-screening. *Journal of Network and Computer Applications* **30**(1) (2007) 99–113
43. Hoffmann, F.: Combining boosting and evolutionary algorithms for learning of fuzzy classification rules. *Fuzzy Sets and Systems* **141** (2004) 47–58
44. Wang, C.H., Hong, T.P., Tseng, S.S.: Integrating membership functions and fuzzy rule sets from multiple knowledge sources. *Fuzzy Sets and Systems* **112** (2000) 141–154
45. Er, M.J., Zhou, Y.: Automatic generation of fuzzy inference systems via unsupervised learning. *Neural Networks* **21**(10) (2008) 1556–1566
46. Chen, C.H., Tseng, V.S., Hong, T.P.: Cluster-based evaluation in fuzzy-genetic data mining. *IEEE Transactions on Fuzzy Systems* **16**(1) (2008) 249–262

47. Nakashima, T., Schaefer, G., Yokota, Y., Ishibuchi, H.: A weighted fuzzy classifier and its application to image processing tasks. *Fuzzy Sets and Systems* **158** (2007) 284–294
48. Ishibuchi, H., Nojima, Y.: Analysis of interpretability-accuracy tradeoff of fuzzy systems by multiobjective fuzzy genetics-based machine learning. *International Journal of Approximate Reasoning* **44**(1) (2007) 4–31
49. Ishibuchi, H., Nakashima, T., Murata, T.: Three-objective genetics-based machine learning for linguistic rule extraction. *Information Sciences* **136** (2001) 109–133
50. Ishibuchi, H., Nakashima, T., Morisawa, T.: Voting in fuzzy rule-based systems for pattern classification problems. *Fuzzy Sets and Systems* **103** (1999) 223–238
51. Lei, Z., Ren-hou, L.: Designing of classifiers based on immune principles and fuzzy rules. *Information Sciences* **178** (2008) 1836–1847
52. Mohamadi, H., Habibi, J., Abadeh, M.S., Saadi, H.: Data mining with a simulated annealing based fuzzy classification system. *Pattern Recognition* **41**(5) (2008) 1824–1833
53. Lowe, D.G.: Similarity metric learning for a variable-kernel classifier. *Neural Computation* **7** (1995) 72–85
54. Chang, C.L.: Finding prototypes for nearest neighbour classifiers. *IEEE Trans Computers* **23** (1974) 1179–1184
55. Tomek, I.: An experiment with the edited nearest-neighbour rule. *IEEE Trans Systems, Man and Cybernetics* **6** (1976) 448–452
56. Tomek, I.: Two modifications of CNN. *IEEE Trans Systems, Man and Cybernetics* **6** (1976) 769–772
57. Verikas, A., Bacauskiene, M., Malmqvist, K.: Learning an adaptive dissimilarity measure for nearest neighbour classification. *Neural Computing & Applications* **11**(3-4) (2003) 203–209
58. Mamdani, E.H., Assilian, S.: An experiment in linguistic synthesis with a fuzzy logic controller. *International Journal of Man-Machine Studies* **7**(1) (1975) 1–12
59. Lin, Y., Cunningham, G.A.: A new approach to fuzzy-neural system modeling. *IEEE Trans Fuzzy Systems* **3**(2) (1995) 190–198
60. Zolghadri, M.J., Mansoori, E.G.: Weighting fuzzy classification rules using receiver operating characteristics (ROC) analysis. *Information Sciences* **177** (2007) 2296–2307
61. Pakkanen, J., Iivarinen, J., Oja, E.: The evolving tree—A novel self-organizing network for data analysis. *Neural Processing Letters* **20** (2004) 199–211
62. Kohonen, T.: *Self-Organizing Maps*. 3 edn. Springer-Verlag, Berlin (2001)
63. Guzaitis, J., Verikas, A.: An efficient technique to detect visual defects in particleboards. *Informatica* **19**(3) (2008) 363–376
64. Song, H.H., Lee, S.W.: LVQ combined with simulated annealing for optimal design of large-set reference patterns. *Neural Networks* **9**(2) (1996) 329–336
65. Kalina, A., Mezyk, E.: Accuracy boosting induction of fuzzy rules with artificial immune systems. In: *Proceedings of the International Multiconference on Computer Science and Information Technology* pp. Volume 3., IEEE (2008) 155–159
66. Halavati, R., Shouraki, S.B., Lotfi, S., Esfandiar, P.: Symbiotic evolution of rule based classifier systems. *International Journal on Artificial Intelligence Tools* **18**(1) (2009) 1–16
67. Ho, S.Y., Chen, H.M., Ho, S.J.: Design of accurate classifiers with a compact fuzzy-rule base using an evolutionary scatter partition of feature space. *IEEE Transactions on Systems Man and Cybernetics, Part B* **34**(2) (2004) 1031–1043