

Fuzzy Sets and Systems 130 (2002) 101-108



www.elsevier.com/locate/fss

A fuzzy neural network for pattern classification and feature selection

Rui-Ping Li^a, Masao Mukaidono^{b, *}, I. Burhan Turksen^c

^aCaelum Research Corporation, Rockville, MD 20850, USA

^bDepartment of Computer Science, Meiji University, 1-1-1 Higash-mita, Tama-ku, Kawasaki-shi 214-8571, Japan ^cDepartment of Industrial Engineering, University of Toronto, Toronto, Canada M5S 1A4

Received 13 August 1997; received in revised form 7 November 2001; accepted 9 January 2002

Abstract

A fuzzy neural network with *memory* connections for classification, and *weight* connections for selection is introduced, thereby solving simultaneously two major problems in pattern recognition: pattern classification and feature selection. The proposed network attempts to select important features from among the originally given plausible features, while maintaining the maximum recognition rate. The resulting value of weight connection represents the degree of importance of feature. Moreover, the knowledge acquired by the network can be described as a set of interpretable rules. The effectiveness of this new method has been validated by using Anderson's IRIS data. The results are: first, the use of two features selected by our method from among the original four in the proposed network results in virtually identical classifier performance; and second, the constructed classifier is described by three simple rules that are of if–then form. © 2002 Elsevier Science B.V. All rights reserved.

Keywords: Pattern classification; Feature selection; Neural networks; Fuzzy logic; Linguistic modeling

1. Introduction

The recognition of patterns is the basis of all science. The aim is to discover a structure in a system consisting of partial subsystems. Usually, something structured refers to the knowledge of the state of a partial subsystem allowing us to easily guess the state of other parts of the same whole system [25]. Techniques of pattern recognition can be generally described as deterministic, statistical, or fuzzy in terms of their axiomatic bases. Traditional

* Corresponding author. Tel: +81-44-934-7450; fax: +81-44-934-7912.

statistical classification methods usually try to find a clear cut boundary to divide the pattern space into some classification regions based on some pre-defined criterion, such as maximizing deviation-betweengroups divided by deviation-within-groups in the linear discriminant analysis (LDA) [10]. As pointed in [2], it is impossible to provide information of degree of uncertainty for a particular example for LDA method since the error rate estimate is a statistical result of the entire sample set. In fact, pattern recognition systems are systems that automatically identify objects based on their measured properties or features derived from these properties. With this viewpoint, a neural network also is a pattern recognition system. The existing neural networks that can be served as classifiers may be grouped into four

E-mail address: masao@cs.meiji.ac.jp (M. Mukaidono).

^{0165-0114/02/\$ -} see front matter 2002 Elsevier Science B.V. All rights reserved. PII: S0165-0114(02)00050-7

categories or their variations: backpropagation (BP) [20], adaptive resonance theory (ART) [4], radial basis functions (RBF) [11], and probabilistic neural networks (PNN) [22]. The first three are based on the deterministic axiomatics, and the last one is based on the probabilistic-statistical axiomatics. Although these techniques have been proven to be useful tools for pattern classification, the selection of features still is a challenge.

Since fuzzy set theory was suggested in 1965 [26], pattern recognition problems have been intensively studied with fuzzy set [3]. The revolutionary significance of fuzzy set theory is that it provides a mathematical method for describing intuitive knowledge of humans. In principle, a mathematical model constructed in accordance with the classical theory must be interpreted in natural language that could be understood intuitively. In contrast to classical methodology, a fuzzy approach to modeling begins with a practical interpretation of concepts, and then generates intuitive logical relations between concepts and constructs a model. A model constructed in accordance with the fuzzy theory, therefore, is certainly interpretable. This methodology is called 'empirical-semantic' approach in [24], and this modeling method is called 'linguistic' modeling in [23]. In recent years, a great deal of attention has been directed toward using the fusion of fuzzy logic and neural networks to develop intelligent systems. This is because the two technologies are strongly complementary to each other [8,19]. Keller [6], for example, incorporated fuzzy membership functions into the Perceptron learning algorithm. Archer [2] used fuzzy set representation in neural network classification problems.

There are two main aspects to the effort of pattern recognition: pattern classification and feature selection. Although many efforts have been made, we still do not have a complete and satisfactory technique that can simultaneously deal with the above two problems. Bezdek [3] proposed a measure of feature selection that works only for binary data. Kuncheva [9] proposed a new selection criterion using the concept of fuzzy rough sets. The latter overcome the limitation of the former; however, combinatorial explosion would become a major problem for the cases in which there are more than a small number of features.

On the other hand, in order to solve the initialization problem and the normalization problem with

traditional learning vector quantization (LVQ) [7], a proportional learning vector quantization (PLVQ) method was introduced in [12,14]. PLVO is a generalized learning vector quantization based on a fuzzy learning law (FLL). Section 2 of this article provides an overview of the PLVO algorithms, since the FLL is employed in the presented network. Section 3 of this article describes our feature-weighted detector (FWD) network that can fulfill both tasks of feature selection and pattern classification. Section 4 includes two examples to verify the effectiveness of our FWD. For the sake of understanding, an artificial data set is chosen in the first example. In the second example, the data set used is Anderson's IRIS data [1] that has been widely studied, allowing us to easily analyze and compare our new technique with existing methods. Finally, Section 5 contains a summary and conclusions.

2. Proportional learning vector quantization

As well known to all, LVQ is a clustering algorithm for organizing a large number of unlabeled vectors into some given clusters. Although some good practical results have been obtained with it, the method still suffers from an initialization problem [18] and a normalization problem [14].

Based on Hebb's learning postulate, we assume that a desired learning rule for weights of LVQ network should satisfy the following differential equation in continuous space:

$$\frac{\mathrm{d}\mathbf{m}_i}{\mathrm{d}t} = \alpha_t u_i(\mathbf{x})(\mathbf{x} - \mathbf{m}_i) \tag{1a}$$

or, in discrete domains, we have

$$\Delta \mathbf{m}_i = \alpha_t u_i(\mathbf{x})(\mathbf{x} - \mathbf{m}_i), \tag{1b}$$

where **x** denotes the input vector, **m**_{*i*} denotes the memory vector of neuron *i*. $u_i(\mathbf{x})$ represents the output value of neuron *i* when **x** is presented in input layer. $\alpha_t = \alpha(1 - t/T)$ is referred to as *temporal* learning rate in [18]. To find the mathematical expression of $u_i(\mathbf{x})$ and its physical meaning, consider the following loss function *L* as introduced in [13,15]:

$$L = \sum_{i=1}^{c} \sum_{k=1}^{N} u_{ik} \|\mathbf{x}_{k} - \mathbf{m}_{i}\|^{2}, \qquad (2)$$

where $u_{ik} \equiv u_i(\mathbf{x}_k)$ represents the degree to which input \mathbf{x}_k ($\mathbf{x}_k = (x_{k1}, x_{k2}, \dots, x_{kp})$) matches memory vector \mathbf{m}_i ($\mathbf{m}_i = (m_{i1}, m_{i2}, \dots, m_{ip})$). *N* is the number of data and *c* is the number of clusters. The number of clusters here is equal to the number of output layer neurons. *p* is the number of features, i.e., the number of input layer nodes. Using the method of the maximum-fuzzy-entropy interpretation [15] and the normalization condition ($\sum_{i=1}^{c} u_{ik} = 1$ for each *k*), the following solution that minimizes the loss function *L* was found:

$$u_{ik} = \exp\left[-\frac{\|\mathbf{x}_k - \mathbf{m}_i\|^2}{2\sigma^2}\right] / \sum_{j=1}^c \exp\left[-\frac{\|\mathbf{x}_k - \mathbf{m}_j\|^2}{2\sigma^2}\right],$$
(3)

where σ is a nonzero number that can be chosen by user. Physically, σ represents the fuzziness in clustering. The smaller the σ is, the less the fuzziness is. There is no theoretical basis for choosing an optimal σ . A heuristic guideline is $\sigma = 0.25\sqrt{p}$ but not limited to this. Note that *p* represents the number of features here. $u_i(\mathbf{x}_k) \in [0, 1]$ is a fuzzy membership function. For a given σ , if the closer the input \mathbf{x}_k is to the memory \mathbf{m}_i , then the closer the output $u_i(\mathbf{x}_k)$ is to one; if the more the input \mathbf{x}_k is away from the memory \mathbf{m}_i , then the closer the output $u_i(\mathbf{x}_k)$ is to zero. From Eqs. (1b) and (3), it is clear that each input updates all the weights (i.e., memory connections $\{\mathbf{m}_i\}$) in proportion as their output values. Eq. (1) was called FLL in [14], and $u_i(\mathbf{x}_k)$ also can be called *special* learning rate corresponding with *temporal learning* rate α_t . When $\sigma \rightarrow 0$, $u_i(\mathbf{x}_k) = \{0, 1\}$, and thus the FLL reduces to competitive learning law (CLL) [7]. The corresponding algorithm is referred to as PLVQ. It has been shown that PLVQ avoids above two problems with LVQ.

The PLVQ Algorithm

- (1) Fix $2 \le c \le n$, $\varepsilon > 0$, $\sigma > 0$ and the maximum number of iterations *T*.
- (2) Initialize $\{\mathbf{m}_i(0)\}$ and learning rate $\alpha_0 \in [0, 1]$.
- (3) For t = 1, 2, ..., T; For k = 1, 2, ..., N;
 (a) Calculate {u_i(x_k)} using Eq. (3).

(b) Update $\{\mathbf{m}_i(t)\}$ based on Eq. (1b), i.e.,

$$\mathbf{m}_{i}(t) = \mathbf{m}_{i}(t-1) + \alpha_{0}(1-t/T)u_{i}(\mathbf{x}_{k})(\mathbf{x}_{k}-\mathbf{m}_{i}).$$
(4)

(c) Next k.

(4) Calculate
$$E = \sum_{i=1}^{c} \sum_{j=1}^{p} |m_{ij}(t) - m_{ij}(t-1)|.$$

(5) IF $E < \varepsilon$ or t > T stop; ELSE next t.

3. Feature-weighted detector networks

A feature-weighted detector (FWD) network is shown in Fig. 1. The network consists of input (I), matching (M), detecting (D) and output (O) layers (Fig. 1(b)). Below, we give a description of this network in detail.

3.1. Input–output relations

As shown in Fig. 1(b), each M node can receive inputs from two sources: the left-right input; and rightleft input from a node of D via a D-M adaptive connection. f() is a comparative function, and the output is the difference of two input values. In detecting layer, there are two types of node: forward and backward nodes. Each forward-node receives p inputs from p nodes of M via pathways with weight connections $\{w_{ii}\}$. g() is Gaussian functions. Each backward-node receives an input from a node of Ovia a backward-pathway with 1 connection fixed. b()is a linear function. The functional of the output layer nodes is to give final classification score for each input by normalizing output values of all D nodes. Each O node receives c+1 inputs. One of them is called set signal. The other c's are from D via c pathways with 1 connection fixed. Set signal occurs before input is presented to the input layer. The role of the set signal is to provide an equal opportunity to match the input for each of M nodes. Before input \mathbf{x}_k is coming, set signal $s_i = 1$, and thus $v_i = s_i = 1$ $(i = 1, 2, \dots, c)$, since b() is a linear function. This guarantees that each of neurons have an equal opportunity to match coming input. When input \mathbf{x}_k is coming, outputs of node of



Fig. 1. (a) A schematic diagram of the FWD network model. (b) Structure and interconnection of neuron i.

neuron *i* are:

$$y_{ij} = (x_{kj} - m_{ji}), \quad j = 1, 2, \dots, p,$$
 (5)

$$z_{i} = \exp\left[-\frac{1}{2\sigma^{2}}\sum_{j=1}^{p}w_{ij}^{2}(x_{kj}-m_{ji})^{2}\right],$$
 (6)

$$u_i = z_i / \sum_{j=1}^{c} z_j.$$
 (7)

3.2. Learning laws

In FWD networks, there are two types of learning when input is presented to the input layer. One is *memory* learning. The other is *weight* learning. \mathbf{m}_i represents the memory of neuron *i*. Memory learning is unsupervised, and the updating rule is based on the FLL, i.e.,

$$\Delta \mathbf{m}_i = \alpha_t u_i(\mathbf{x}_k) (\mathbf{x}_k - \mathbf{m}_i), \qquad (8)$$

where \mathbf{x}_k represents the *k*th input. On the other hand, in weight learning w_{ij} represents the degree to which feature *j* contributes to the cluster *i*. In order to find the updating rule of $\{w_{ij}\}$, introduce the following error function:

$$E = \frac{1}{2} \sum_{k=1}^{N} \sum_{i=1}^{c} (u_i(\mathbf{x}_k) - d_i)^2,$$
(9)

where d_i is the desired value of output layer node *i*. Therefore, unlike memory learning, weight learning is supervised. Based on the chain rule of differential calculus, using Eqs. (9), (7) and (6) the following updating rule is obtained:

$$\Delta w_{ij} = \frac{\beta}{\sigma^2 s^2} \left(u_i(\mathbf{x}_k) - d_i \right)$$
$$\times t \left(\sum_{j=1}^c z_j - z_i \right) w_{ij} z_i (x_{kj} - m_{ji})^2, \qquad (10)$$

where $\beta > 0$ is learning rate. For the sake of understanding, design $0 \le w_{ij} \le 1$ for each *i* and *j*. $w_{ij} = 0$ means that feature *j* has no 'contribution' to cluster *i*; and $w_{ij} = 1$ means that feature *j* has the most contribution to cluster *i*. The algorithm can be stated as follows.

Feature-weighted Detector (FWD) Network Algorithm

- 1. Fix $\sigma > 0$, $\alpha \in [0, 1]$, $\beta > 0$, $\varepsilon > 0$, and the maximum number of iterations *T*.
- 2. Initialize $\{\mathbf{m}_i(0)\}$, using *c* samples randomly chosen from $\{\mathbf{x}_k\}$ (k = 1, 2, ..., N), and $w_{ij}(0) = 1$ for each *i* and *j*.
- 3. For t = 1, 2, ..., T; For k = 1, 2, ..., N
 - (a) Calculate $\{u_i\}$ using Eq. (7).
 - (b) Update $\{\mathbf{m}_i(t)\}$ using Eq. (8).





Fig. 2. Data set of Example 1.

- (c) Update $\{w_{ij}(t)\}$ using Eq. (10). (d) Next k.
- 4. Calculate *E* using Eq. (9).
- 5. IF $E < \varepsilon$, or t > T stop, ELSE next t.

4. Applications

Two examples have been selected to illustrate the performance of our FWD network. The purpose of selecting Example 1 is to help us to intuitively understand the physical meaning of the method. From the result of Example 2, a potential application value of the presented method should be shown.

4.1. Example 1—an artificial data set

For simplicity and intuition, we applied an artificial data set in which each pattern has two plausible features as showed in Fig. 2 to FWD. From the distribution of the data, it is clear that feature x_1 of this example has no contribution to the classification that follows the target outputs of Table 1.

For data set of Fig. 2, we run the FWD with $\varepsilon = 0.058$, $\alpha = 0.01$, $\beta = 0.1$ and $\sigma = 0.35$. In this example, the number of input layer nodes is two (i.e., p = 2), and the number of output layer nodes is also two (i.e., c = 2). After 218 iterations, for each input the actual output is listed in the right two columns of Table 1. This classification is fuzzy. The item 'target output' of Table 1 list the desired output when input k is presented to the input layer of the FWD. In the Table 1, u_i represents the output of neuron i. Here one neuron corresponds to one cluster. Each pattern

should be assigned to either of two clusters in hard classification as shown in the item of 'target output' of Table 1. Obviously, if transforming the actual output listed in the right two columns of Table 1 into 0-1 binary value, then the actual classification result of the FWD is identical with the target in this example. Preferably, after learning, the resulting weight connections of the FWD are $\mathbf{w}_1 = (0.10, 0.99)$ and $\mathbf{w}_2 = (0.15, 0.99)$. The contribution of feature x_1 to the cluster 1 is $w_{11} = 0.10$; the contribution of feature x_1 to the cluster 2 is $w_{21} = 0.15$; the contribution of feature x_2 to the cluster 1 is the same as that to the cluster 2, i.e., $w_{12} = w_{22} = 0.99$. It shows that feature x_1 can be eliminated from the plausible features set selected initially since the degree of the importance of feature x_1 is much less than that of feature x_2 . Therefore, the presented FWD enables the classification of pattern, and as well the selection of feature.

4.2. Example 2—an application to IRIS data

IRIS data [1] has been used in many papers to illustrate various clustering methods [21,17]. The motivation of selecting IRIS data here is since we have already known the typical performance of the existing methods applied to it and also we can analyze feature by means of the geometric structure as used in [3]. As well known, the IRIS data are a set of 150 fourdimensional vectors. The plausible features selected initially include sepal length (x_1) , sepal width (x_2) , petal length (x_3) and petal width (x_4) . The 150 IRIS data used come from three subspecies (clusters): sestosa, versicolor, and virginica. Each subspecies owns 50 samples, respectively. Anderson measured each feature of 150 plants (samples). For this data, using the existing methods the typical number of mistakes is around 5 for supervised classifiers, and around 15 for unsupervised classifiers [17].

Shown in Table 2 are results of two experiments. In the first experiment (Experiment 1), all of four plausible features were used. The results are: the number of mistakes, e = 5; three weight vectors are $\mathbf{w}_1 =$ (1.00, 1.00, 0.95, 1.00), $\mathbf{w}_2 = (0.00, 0.00, 1.00, 1.00)$, and $\mathbf{w}_3 = (0.00, 0.00, 0.82, 0.97)$. Based on this result, it has been clearly shown that (1) features x_1 and x_2 have no contribution to clusters s_2 and s_3 , and (2) the role of features x_1 and x_2 in cluster s_1 also can be jointly played by features x_3 and x_4 . For this,

Table 1							
Classification	results	of	data	set	of	Fig.	2

Data		Target outpu	Target outputs		Outputs of neurons		
k	X	$\overline{u_1}$	<i>u</i> ₂	$\overline{u_1}$	<i>u</i> ₂		
1	(0.05, 0.12)	1.00	0.00	0.989	0.011		
2	(0.10, 0.20)	1.00	0.00	0.901	0.099		
3	(0.15, 0.05)	1.00	0.00	0.998	0.002		
4	(0.20, 0.10)	1.00	0.00	0.993	0.007		
5	(0.20, 0.20)	1.00	0.00	0.899	0.101		
6	(0.70, 0.35)	0.00	1.00	0.118	0.882		
7	(0.75, 0.45)	0.00	1.00	0.008	0.992		
8	(0.80, 0.40)	0.00	1.00	0.033	0.967		
9	(0.85, 0.10)	0.00	1.00	0.002	0.998		
10	(0.80, 0.10)	1.00	0.00	0.993	0.007		
11	(0.83, 0.15)	1.00	0.00	0.974	0.026		
12	(0.85, 0.19)	1.00	0.00	0.925	0.075		
13	(0.90, 0.07)	1.00	0.00	0.997	0.003		
14	(0.87, 0.13)	1.00	0.00	0.985	0.015		
15	(0.10, 0.35)	0.00	1.00	0.122	0.872		
16	(0.15, 0.45)	0.00	1.00	0.008	0.992		
17	(0.20, 0.40)	0.00	1.00	0.032	0.968		
18	(0.25, 0.50)	0.00	1.00	0.002	0.998		

Table 2 Experimental results for IRIS data set

	Experiment 1 $\{x_1, x_2, x_3, x_4\}$	Experiment 2 $\{x_3, x_4\}$
σ	0.50	0.50
α	0.01	0.01
β	0.10	0.10
3	2.00	2.00
Т	1000	1000
е	5	5
\mathbf{w}_1	(1.00, 1.00, 0.95, 1.00)	(1.00, 1.00)
w ₂	(0.00, 0.00, 1.00, 1.00)	(1.00, 1.00)
w ₃	(0.00, 0.00, 0.82, 0.97)	(0.82, 0.97)

Features x_1 , x_2 , x_3 , and x_4 represent sepal length, sepal width, petal length, and petal width, respectively. *e* represents the number of mistakes in classification, and $\{\mathbf{w}_i\}$ represent weight vectors.

features x_1 and x_2 is supposed being meaningless. In order to prove that, in the second experiment (Experiment 2) only features x_3 and x_4 were used. The results of the second experiment are demonstrated in the right column of Table 2: the number of mistakes of the second experiment is the same as that of the first experiment, i.e., e = 5; three weight vectors are $\mathbf{w}_1 = (1.00, 1.00), \mathbf{w}_2 = (1.00, 1.00), \mathbf{w}_3 = (0.82, 0.97).$

Obviously, the use of two features x_3 and x_4 results in virtually identical classifier performance. After feature selection, therefore, only feature x_3 (petal length) and feature x_4 (petal width) are chosen.

For the sake of description, above-selected two feature variables are renamed: x_1 now represents petal length, and x_2 petal width. When a new input $\mathbf{x} = (x_1, x_2)$ is coming, using the obtained $\{w_{ij}\}, \{m_{ji}\}$

106

and given σ , based on Eqs. (6) and (7), we have the following:

$$z_{i} = \exp\left[-\frac{1}{2\sigma^{2}} w_{i1}^{2} (x_{1} - m_{1i})^{2}\right]$$

$$\times \exp\left[-\frac{1}{2\sigma^{2}} w_{i2}^{2} (x_{2} - m_{2i})^{2}\right]$$

$$= U_{A_{i1}}(x_{1}) \times U_{A_{i2}}(x_{2}), \quad i = 1, 2, 3, \quad (11)$$

$$u_{c_i} = U_{A_{i1}}(x_1)U_{A_{i2}}(x_2) \sum_{j=1}^3 U_{A_{j1}(x_1)U_{A_{j2}}(x_2)},$$

$$i = 1, 2, 3.$$
(12)

Note that z_1 , z_2 and z_3 are outputs that correspond, respectively, *class* sestosa (c_1) , *class* versicolor (c_2) , and *class* virginica (c_3) . Normalized $u_{c_i}(\mathbf{x})$ represents the degree to which input \mathbf{x} belongs to *class* c_i (i = 1, 2, 3). Further, we have

$$U_{A_{11}}(x_1) = \exp[-2(x_1 - 1.46)^2], \tag{13}$$

$$U_{A_{12}}(x_2) = \exp[-2(x_2 - 0.25)^2], \qquad (14)$$

$$U_{A_{21}}(x_1) = \exp[-2(x_1 - 4.29)^2], \tag{15}$$

$$U_{A_{22}}(x_2) = \exp[-2(x_2 - 1.36)^2],$$
(16)

$$U_{A_{31}}(x_1) = \exp[-1.34(x_1 - 5.54)^2], \tag{17}$$

$$U_{A_{32}}(x_2) = \exp[-1.88(x_2 - 2.0)^2].$$
(18)

The following three fuzzy rules, which correspond, respectively, to i = 1, 2, 3 in Eq. (11), are thus constructed:

- R_1 IF petal length is *nearly* 1.46 and petal width is *nearly* 0.25, THEN it is *sestosa*.
- R_2 IF petal length is *nearly* 4.29 and petal width is *nearly* 1.36 THEN it is *versicolor*.
- R_3 IF petal length is *nearly* 5.54 and petal width is *nearly* 2.00, THEN it is *virginica*.

Note that IF-part in above rules corresponds to the right-hand side of Eq. (11), and THEN-part the left-hand side of Eq. (11), where fuzzy numbers, *nearly* 1.46, *nearly* 0.25, *nearly* 4.29, *nearly* 1.36, *nearly* 5.54, and *nearly* 2.00 are, respectively, represented by fuzzy sets A_{11} , A_{12} , A_{21} , A_{22} , A_{31} , and A_{33} in Eq. (11). The numbers, 1.46, 0.25, 4.29, 1.36, 5.54, and 2.00 are derived from Eqs. (13)–(18).

5. Conclusions

A fuzzy neural network that enables the classification of patterns and the selection of features is introduced. This algorithm includes two types of learning, i.e., unsupervised learning for memory connection and supervised learning for weight connection. Examples that are provided has demonstrated the ability of the feature-weighted detector (FWD) network to classify pattern and select feature. Moreover, distinct to traditional neural networks for which it is usually difficult to interpret the obtained knowledge [5,16], our FWD provides interpretable rules that are of if-then form. These properties of the FWD suggest that it will be a promising method for pattern recognition.

References

- E. Anderson, The IRISes of the Gaspe Peninsula, Bull. Amer. IRIS Soc. 59 (1939) 2–5.
- [2] N.P. Archer, S. Wang, Fuzzy set representation of neural network classification boundaries, IEEE Trans. Systems Man Cybernet. 21 (1991) 735–742.
- [3] J.C. Bezdek, Pattern Recognition with Fuzzy Objective Function Algorithms, Plenum Press, New York, 1981.
- [4] G.A. Carpenter, S. Grossberg, Pattern Recognition by Self-organizing Neural Networks, MIT Press, Cambridge, MA, 1991.
- [5] L. Fu, Rule generation from neural networks, IEEE Trans. Systems Man. Cybernet. 24 (8) (1994) 1114–1124.
- [6] J.M. Keller, D.T. Hunt, Incorporating fuzzy membership functions into the perceptron algorithm, IEEE Trans. PAMI PAMI-7 (6) (1985) 693–699.
- [7] T. Kohonen, Self-organization and Associative Memory, Springer, Berlin, 1989.
- [8] B. Kosko, Neural Networks and Fuzzy Systems, Prentice-Hall, Englewood Cliffs, NJ, 1991.
- [9] L.I. Kuncheva, Fuzzy rough sets: application to feature selection, Fuzzy Sets and Systems 51 (1992) 147–153.
- [10] P.A. Lachenbruch, Discriminant Analysis, Hafner, New York, 1975.
- [11] S. Lee, R.M. Kil, A Gaussian potential function network with hierarchically self-organizing learning, Neural Networks 4 (2) (1991) 207–224.
- [12] R.-P. Li, M. Mukaidono, Proportional learning law and local minimum escape in clustering networks, International Conference on Neural Information Processing ICONIP'95, Beijing, 1995, pp. 684–686.
- [13] R.-P. Li, M. Mukaidono, A maximum-entropy approach to fuzzy clustering, International Joint Conference of FUZZ-IEEE/IFES'95, 1995, pp. 2227–2233.
- [14] R.-P. Li, M. Mukaidono, Proportional learning vector quantization, Journal of Japan Society for Fuzzy Theory and Systems 10 (6) (1998) 1129–1134.

- [15] R.-P. Li, M. Mukaidono, Gaussian clustering method based on maximum-fuzzy entropy interpretation, Fuzzy Sets and Systems 102 (1999) 153–258.
- [16] S. Mitra, Y. Hayashi, Neuro-fuzzy rule generation: survey in soft computing framework, IEEE Trans. Neural Networks 11
 (3) (2000) 748–768.
- [17] S.C. Newton, S. Pemmaraju, S. Mitra, Adaptive fuzzy leader clustering of complex data sets in patterns recognition, IEEE Trans. Neural Networks 3 (5) (1992) 794–800.
- [18] N.R. Pal, J.C. Bezdek, E.C.-K. Tsao, Generalized clustering networks and Kohonen's self-organizing scheme, IEEE Trans. Neural Networks 4 (1993) 549–557.
- [19] S.K. Pal, S. Mitra, Neuro-Fuzzy Pattern Recognition: Methods in Soft Computing, Wiley, New York, 1999.
- [20] D.E. Rumelhart, J.L. McCleland, Parallel Distributed Processing, MIT Press, Cambridge, MA, 1986.

- [21] P.K. Simpson, Fuzzy min-max neural networks-part 1: classification, IEEE Trans. Neural Networks 3 (5) (1992) 776–786.
- [22] D.F. Specht, Probabilistic neural networks, Neural Network 3 (1990) 109–118.
- [23] M. Sugeno, T. Yasukawa, A fuzzy-logic-based approach to qualitative modeling, IEEE Trans. Fuzzy Systems 1 (1993) 7–31.
- [24] I.B. Turksen, Measurement of membership functions and their acquisition, Fuzzy Sets and Systems 40 (1991) 5–38.
- [25] S. Watanabe, Pattern Recognition, Wiley Interscience, New York, 1985.
- [26] L.A. Zadeh, Fuzzy Sets, Inform. and Control 8 (1965) 338–353.