

Learning of Class Membership Values by Ellipsoidal Decision Regions

Leechter Yao and Chin-Chin Lin

Abstract—A novel method of learning complex fuzzy decision regions in the n -dimensional feature space is proposed. Through the fuzzy decision regions, a given pattern's class membership value of every class is determined instead of the conventional crisp class the pattern belongs to. The n -dimensional fuzzy decision region is approximated by union of hyperellipsoids. By explicitly parameterizing these hyperellipsoids, the decision regions are determined by estimating the parameters of each hyperellipsoid. Genetic Algorithm is applied to estimate the parameters of each region component. With the global optimization ability of GA, the learned decision region can be arbitrarily complex.

Keywords— ellipsoid, genetic algorithm, decision regions, classification

I. INTRODUCTION

PATTERN classification mainly deals with determination of decision regions based on the given prototypes. The information carried by the every prototype consists of the features associated with the prototype and the class the prototype belongs to. Since the information carried by each prototype are gathered by human beings, it is understood that uncertainty might exist within the information assigned to the prototype. Hence, fuzziness could be involved in feature space or in class assignment. For the ease of analysis and manipulation, most of the research in the field [1-7] gives crisp feature descriptions yet leave class assignment or classification fuzzily defined. Different from the deterministic classification [8-9] where each prototype is classified into one and only one class, membership degrees are employed in fuzzy classification defining degrees of belonging of each prototype to every class. There have been numerous approaches proposed for clustering [1, 10-14] based on the prototypes with fractional membership degrees belonging to different classes. However, not too many researches investigate learning of fuzzy decision regions based on the prototypes with fractional membership degrees belonging to every different class. The conventional fuzzy classification approach assigns each prototype's degree of belonging to different class by a real number between 0 and 1. The real number is called the class membership value. The

training of conventional fuzzy classification approach is to find the model that determines the crisp classification based on the training prototype. For some classification applications such as medical diagnosis, geographical analysis or decision making, the classification aims to generate each pattern's class membership values for each class other than the crisp pattern class. For some applications, the crisp class that a given pattern is classified may not be as important as the class membership values since the class membership values will serve as important decision support data. In the paper, a novel approach is proposed to the train the decision region so that the class membership value of every class can be determined through the decision regions. A set of hyperellipsoids with adaptively tuned centers, orientations and sizes are employed to learn the fuzzy decision regions. If the prototypes are classified into c classes, there are in fact c decision regions to be learned since each prototype has crisply defined coordinates in the feature space and c respectively assigned membership degrees belonging to each of c classes. The proposed learning scheme for decision regions is basically nonparametric since no statistical information of the prototypes is assumed. The decision regions determined by the perceptron algorithm [15-16] or the least mean square algorithm [8-9] are compose of half spaces. Hyperellipsoids geometrically locates and cover the decision regions more precisely and yet requires more concise parameterizations.[17-18]

Since the decision region can be approximated by an union of a finite number of hyperellipsoids, learning of complex decision region is equivalent to the estimation of the parameters of hyperellipsoids. Multiple hyperellipsoids are respectively employed to approximate the decision regions in [19] and [20]. In [19], a multivariate Gaussian distribution function of prototypes is assumed. The locus of prototypes with constant probability density function forms a hyperellipsoid in the feature space. Therefore, [19] aims to learn the distribution function of prototypes, which is equivalent to learning the parameters of hyperellipsoids. The Genetic Algorithm (GA) is utilized to learn the hyperellipsoids in [20]. Both [19] and [20] consider only the prototypes with crisp classification, i.e., each prototype is assigned one and only one class. In this paper, the decision region for the prototypes with fuzzy classification will be investigated.

Similar to [20], the GA is also to be used in this paper to estimate the parameters of ellipsoids. The number of hyperellipsoids required to approximate the decision region is generally unknown in advance. More hyperellipsoids than

Manuscript received Nov. 19, 2004.

Leechter Yao is with Dept. of Electrical Engineering, National Taipei University of Technology, Taipei, Taiwan. (phone : +886-2-2751-0623; fax : +886-2-2751-3892; e-mail : ltyao@ntut.edu.tw).

Chin-Chin Lin is with Dept. of Electrical Engineering, National Taipei University of Technology, Taipei, Taiwan. (e-mail : s0669013@ntut.edu.tw).

learning fuzzy decision region is equivalent to tuning the parameters of region components in Q_i so that

$$Q_i \equiv \bigcup_{p=1}^{b_i} (S_{ip} \cap W) \approx G, i = 1 \dots c_i \quad (11)$$

III. LEARNING OF FUZZY DECISION REGIONS

A. Fitness Function of GA

GA learns different sets of region components Q_1, Q_2, \dots, Q_c to approximate the fuzzy decision regions G_1, G_2, \dots, G_c . The approximation is performed in the sense of minimum misclassification errors as well as minimum total volume of region components. GA is implemented to learn the fuzzy decision regions G_1, G_2, \dots, G_c in parallel while satisfying the constraints of fuzzy partitions in (1) and (2) at the same time.

Since the volume of a region component S_{ip} is proportional to the determinant of A_{ip}^{-1} , the fitness function for GA to learn the parameters of region components approximating G_i is given as

$$e = e_a + \gamma \sum_{i=1}^c \sum_{p=1}^{b_i} \prod_{k=1}^n d_{ipk}^2 \quad (12)$$

where e_a denotes the misclassification errors due to the set of region components; and γ weights the total volume of region components with respect to the misclassification errors. Based on the misclassification errors as well as the total volume of region components, GA is able to tune the parameters of each region component to minimize the misclassification errors as well as to adjust the sizes and orientations of region components to geometrically approximate the fuzzy decision region in parallel.

Referring to (3)-(10), the parameters of a region component S_{ip} to be learned by GA are the coordinates of center $\mathbf{v}_{ip} \equiv [v_{ip1}, v_{ip2}, \dots, v_{ipn}]^T$, the angles of rotations in each direction $\boldsymbol{\theta}_{ip} \equiv [\theta_{ip1}, \theta_{ip2}, \dots, \theta_{ipn}]^T$ and the length of each axes of the ellipsoid $\mathbf{d}_{ip} \equiv [d_{ip1}, d_{ip2}, \dots, d_{ipn}]^T$, $i = 1 \dots c, p = 1 \dots b_i$. In order to learn the fuzzy decision region G_i , it is shown in (11) that more than enough region components S_{ip} are set to learn by GA. The S-norm of an union in (11) is operated by $\max(\cdot)$. Referring to (11) and (3), the distance between the j -th prototype and the region component Q_i can be calculated by:

$$y_{ij} = \bigvee_{p=1}^{b_i} f_{ip}(\mathbf{x}_j; \mathbf{v}_{ip}, \boldsymbol{\theta}_{ip}, \mathbf{d}_{ip}) \quad (13)$$

$$= \bigvee_{p=1}^{b_i} (\mathbf{x}_j - \mathbf{v}_{ip})^T \boldsymbol{\theta}_{ip} A_{ip} (\mathbf{x}_j - \mathbf{v}_{ip})$$

where \bigvee denotes the operation of maximization. Similar to the fuzzy partition techniques utilized in fuzzy c-means algorithm, membership degrees are determined by minimizing the following objective function adjoining the constraint in (1):

$$J = \sum_{i=1}^c \sum_{j=1}^m \hat{u}_{ij}^{\eta_i} y_{ij} + \sum_{j=1}^m \lambda_j (1 - \sum_{i=1}^c \hat{u}_{ij}) \quad (14)$$

where \hat{u}_{ij} is the estimated membership degree associated with y_{ij} in (13), η_i is a weighting exponent which determines the fuzziness of membership degrees \hat{u}_{ij} , and λ_j is the Lagrange

multiplier. Membership degrees \hat{u}_{ij} are determined by minimizing the objective function in (14); they thus can be obtained by setting

$$\frac{\partial J}{\partial \hat{u}_{ij}} = 0, \quad (15)$$

$$\text{and } \frac{\partial J}{\partial \hat{u}_{ij}} = 0, \quad (16)$$

Solving \hat{u}_{ij} in (15) and (16), gives

$$\hat{u}_{ij} = \frac{\left(\frac{1}{y_{ij}}\right)^{\frac{1}{\eta_i-1}}}{\sum_{i=1}^c \left(\frac{1}{y_{ij}}\right)^{\frac{1}{\eta_i-1}}}, i = 1 \dots c, j = 1 \dots m \quad (17)$$

The total misclassification errors e_a in (12) can be calculated by

$$e_a = \sum_{i=1}^c \sum_{j=1}^m (u_{ij} - \hat{u}_{ij})^2 \quad (18)$$

In order to calculate the misclassification error for each prototype (\mathbf{x}_j, u_{ij}) , the distance between \mathbf{x}_j and region components Q_1, Q_2, \dots, Q_c , denoted by $y_{1j}, y_{2j}, \dots, y_{cj}$, are first calculated by (13). The membership degrees of the given prototype belonging to each class are then calculated by (17) based on $y_{1j}, y_{2j}, \dots, y_{cj}$. Finally, the misclassification error corresponding to the given prototype is calculated by (18). Therefore, although region components Q_1, Q_2, \dots, Q_c are learned in parallel to approximate the fuzzy decision regions G_1, G_2, \dots, G_c , the learning for Q_1, Q_2, \dots, Q_c are cross correlated.

Note that since the weighting exponent η_i (≥ 1) has significant influence on the fuzziness of fuzzy partition, it is also learned by GA along with the parameters of each region component, i.e., $\mathbf{v}_{ip}, \boldsymbol{\theta}_{ip}$ and \mathbf{d}_{ip} , $i = 1 \dots c, p = 1 \dots b_i$, based on the fitness function in (12) and (18).

In fuzzy c-means algorithm, cluster center is iteratively adjusted so that the total distance norms between the prototype and every cluster center are minimized. In (13), the Mahalanobis norm is used for measuring the distance between the prototype and the center of Q_i, S_{ip} , $i = 1 \dots c, p = 1 \dots b_i$. Similar to fuzzy c-means algorithm, a virtual cluster center \mathbf{v}'_i and the virtual distance norm D_{ijA}^2 are induced from (13) for each fuzzy decision region G_i corresponding to the prototypes (\mathbf{x}_j, u_{ij}) , $i = 1 \dots c, j = 1 \dots m$, i.e.,

$$D_{ijA}^2 = \|\mathbf{x}_j - \mathbf{v}'_i\|_A^2 = \bigvee_{p=1}^{b_i} (\mathbf{x}_j - \mathbf{v}_{ip})^T \boldsymbol{\theta}_{ip} A_{ip} (\mathbf{x}_j - \mathbf{v}_{ip}) \quad (19)$$

Instead of iteratively determining the cluster center as in fuzzy c-means algorithm, centers as well as other parameters of Q_i are iteratively determined by GA in this paper. Since the distance norm in (19) is not a linear function, centers of Q_i cannot be calculated by differentiating the objective function with respect to the parameters of region component center (\mathbf{v}_{ip}) as in the regular fuzzy c-means algorithm. GA is thus utilized instead as the tool of optimization. Both fuzzy c-means algorithm and the method proposed in this paper, aim to minimize the total distance norms as in (19) between prototypes and each cluster center. For fuzzy c-means algorithm, the minimization of distance norms are performed by iteratively tuning cluster center as the weighted topological mean among the given

prototypes, where the weights are the membership degrees calculated based on the distance norms between prototypes and cluster center. The minimization of distance norms proposed in this paper is nevertheless achieved by iteratively tuning a set of candidate hyperellipsoids to cover the prototypes based on the membership degree associated with each prototype.

B. Trimming of Redundant Region Components

Since suitable number of region components required to learn the fuzzy decisions is generally unknown. To begin with, more than enough region components are assigned to learn the fuzzy decision regions. It might be possible that GA tunes the parameters of region components so that less number of region components are well conglomerated to approximate the fuzzy decision region. In other words, redundant region components might exist after the learning process by GA. It is thus necessary to design a scheme to trim off those redundant region components to increase the classification accuracy. In order to determine the redundant region components, the contribution degree of each estimated region component is defined. Referring to (13), b_i region components are set to conglomerate as Q_i to approximate the fuzzy decision region G_i . In the feature space $W \subset R^n$, there is one and only one region component that is closest to a prototype x_j , $j = 1 \dots m$. For any p -th region component learning the i -th fuzzy decision region, S_{ip} , $p \in [1, b_i]$, define the threshold function $T(\cdot)$ as

$$T(x_j, S_{ip}) = \begin{cases} 1, & \text{if } f_{ij}(x_j; d_{ip}, v_{ip}, \theta_{ip}) > f_{ij}(x_j; d_{iq}, v_{iq}, \theta_{iq}), \forall q \in [1, b_i], q \neq p; \\ 0, & \text{otherwise} \end{cases} \quad (20)$$

The contribution degree of an estimated region component S_{ip} can be defined as the number of the prototypes that are closest to it. Let $g(S_{ip})$ be the contribution degree of S_{ip} , $i = 1 \dots c$, $p = 1 \dots b_i$, then

$$g(S_{ip}) = \sum_{j=1}^m T(x_j, S_{ip}) \quad (21)$$

With the contribution degree given as in (21), the redundant region components can thus be defined as the ones with contribution degree less than a preset value. Within Q_i , let the redundant region components be \bar{S}_{ir} , $r \in [1, b_i]$, then

$$g(\bar{S}_{ir}) < \alpha \quad (22)$$

where α is a preset threshold value determining whether the evaluated region component is considered to be redundant. If β_i region components are determined to be redundant, trimming off these redundant region components from Q_i will thus increase the classification efficiency since fewer coefficients are required to parameterize fuzzy decision regions. Let \hat{Q}_i be the refined set of candidate region components to approximate fuzzy decision region G_i , then referring to (11),

$$\hat{Q}_i = Q_i - \bigcup_{i=1}^{\beta_i} (\hat{S}_{ip} \cap W), \quad i = 1 \dots c \quad (23)$$

IV. IMPLEMENTATION OF THE GENETIC ALGORITHM

Each n -dimensional region component S_{ip} is parameterized by the coordinates of center v_{ip} , the angles of rotations in each direction θ_{ip} and the length of each axes d_{ip} , $i = 1 \dots c$, $p = 1 \dots b_i$.

Along with the weighting exponent η_i , there are $(3n+1)$ parameters to be learned by GA for each region component. Every estimated parameter is encoded as a string of binary digits. The binary strings are then cascaded to form a chromosome. If c decision regions are to be determined from the information carried by prototypes, only $(c-1)$ sets of parameters need to be learned since the c -th decision region can be determined based on the constraint in (1). That is, membership degree of the c -th decision region is determined by subtracting all the membership degrees belonging to other classes from 1. Therefore, total number of parameters encoded in one chromosome is given by

$$\kappa = (3n+1) \sum_{i=1}^{c-1} b_i \quad (24)$$

Referring to (12) and (18), fuzzy decision regions G_1, G_2, \dots, G_c are learned by optimizing v_{ip} , θ_{ip} , d_{ip} and η_i via GA, i.e.

$$(v_{ip}, \theta_{ip}, d_{ip}, \eta_i)_{i=1..c} = \underset{i=1..c, p=1..b_i}{\operatorname{argmin}} \left(\sum_{i=1}^c \sum_{j=1}^m (u_{ij} - \hat{u}_{ij})^2 + \gamma \sum_{i=1}^c \sum_{p=1}^{b_i} \prod_{k=1}^n d_{ipk}^2 \right) \quad (25)$$

V. NUMERICAL SIMULATION

In this section, a numerical example is presented verifying the proposed algorithm. It will be shown that the number of region components required to approximate the decision region for every class are unknown a priori. Although more than enough region components are assigned to learn the decision region, the redundant region components can be easily identified and trimmed off based on the learning results. The weighting coefficient γ in (12) is set to be 0.15. In order to assess the learning efficiency and accuracy, the prototypes are divided into two parts; one part for the learning and the other part for the test. In this example, 695 prototypes with 10% noise are set for the learning. The membership degrees for class 1, 2 and 3 of these 695 prototypes are shown in Fig. 1, respectively. The distribution of prototypes in this example is not linearly separable. The decision region of every class is more complicated and more difficult to learn. In this example, 12 region components are assigned, respectively, to learn each of 3 classes. The learning results that 6 out of 12 candidate region components for both class 1 and 2, 5 out of 12 candidate region components for both class 3 remain based on contribution degree of each region component. The decision regions for class 1, 2 and 3 are shown in Fig. 2, respectively. As the decision regions are learned based on the prototypes, another set of 695 data is employed for test. The average misclassification error \bar{e} is calculated to be 0.00881. The learning of decision regions is still accurate and efficient.

VI. CONCLUSION

It has been shown in this paper that the desired fuzzy decision region is approximated by a finite number of ellipsoids. By appropriately parameterizing the hyperellipsoids, the GA is applied to estimate the associated parameters, and thus to learn

the decision region. Since the minimization criterion of GA is defined to be a combination of misclassification error and the sum of the volume of the estimated hyperellipsoids, the learning hyperellipsoids tends to agglomerate together with the least total volume to approximate the decision regions. Compared to the traditional methods such as statistical approaches or artificial neural networks that approximate decision regions with half spaces, the proposed method locates and approximates the decision regions more precisely and yet employs more concise parameterizations.

REFERENCES

- [1] J. C. Bezdek, *Pattern Recognition with Fuzzy Objective Function Algorithms*, New York: Plenum Press, 1981.
- [2] J. C. Bezdek, S. K. Pal, eds., *Fuzzy Models for Pattern Recognition*, New York: IEEE Press, 1992
- [3] J. C. Bezdek, "Computing with uncertainty," *IEEE Commun. Mag.*, vol. 30, Sept. 1992, pp. 24-37
- [4] S. K. Pal, D. K. Dutta Majumber, *Fuzzy Mathematical Approach to Pattern Recognition*, New York: John Wiley, 1986.
- [5] S. K. Pal, S. Mitra, Multilayer perception, fuzzy sets and classification, *IEEE Trans. Neural Networks*, vol. 3, May 1992, pp. 683-697
- [6] W. Pedrycz, "Fuzzy sets in pattern recognition: methodology and methods," *Pattern Recognition*, vol. 23, 1990, pp. 121-146
- [7] W. Pedrycz, "Fuzzy neural networks with reference neurons as pattern classifiers," *IEEE Trans. Neural Networks*, vol. 3, May 1992, pp. 770-775
- [8] S. T. Bow, *Pattern Recognition Application to Large Dataset Problems*, New York: Marcel Dekker, 1980.
- [9] R. S. Chalkoff, *Pattern Recognition: Statistical, Structural and Neural Approaches*, New York: John Wiley & Sons, 1992.
- [10] J. C. Bezdek, J. M. Keller, R. Krishnapuram, N. R. Pal, *Fuzzy Models and Algorithms for Pattern Recognition and Image Processing*, Norwell, MA: Kluwer, 1999.
- [11] R. J. Hathaway, J. C. Bezdek, Y. Hu, "Generalized fuzzy c-means clustering strategies using L_p norm distances," *IEEE Trans. Fuzzy Syst.*, vol. 8, May 2000, pp. 576-582
- [12] R. Krishnapuram, J. M. Keller, "A possibility approach to clustering," *IEEE Trans. Fuzzy Syst.*, vol. 1, Mar. 1993, pp. 98-110
- [13] L. Zhao, Y. Tsujimura, M. Gen, "Genetic algorithm for fuzzy clustering," *Proceeding of IEEE International Conference on Evolutionary Computation*, 1996, pp. 716-719.
- [14] B. P. Buckles, F. E. Petry, D. Prabhu, R. George and R. Srikanth, "Fuzzy clustering with genetic search," *Proceeding of IEEE World Congress on Computational Intelligence*, 1994, pp. 46-50.
- [15] R. P. Lippmann, "An introduction to computing with neural nets," *IEEE ASSP Mag.*, April 1987, pp. 4-22.
- [16] T. Khanna, *Foundations of Neural Networks*, MA: Addison-Wesley, 1990.
- [17] H. I. Avi-Itzhak, J. A. Van Mieghem, L. Rub, "Multiple subclass pattern recognition: a maximum correlation approach," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 17, April 1995, pp. 418-431.
- [18] Q. Zhu, Y. Cai, "A subclass model for nonlinear pattern classification," *Pattern Recognition Lett.*, vol. 19, Feb. 1998, pp. 19-29.
- [19] Q. Zhu, Y. Cai, L. Liu, "A multiple hyper-ellipsoidal subclass model for an evolutionary classifier," *Pattern Recognition*, vol. 34, March 2001, pp. 547-560.
- [20] L. Yao, "Nonparametric learning of decision regions via the genetic algorithm," *IEEE Trans. System, Man, and Cybernetics*, vol. 26, Feb. 1996, pp. 313-321.

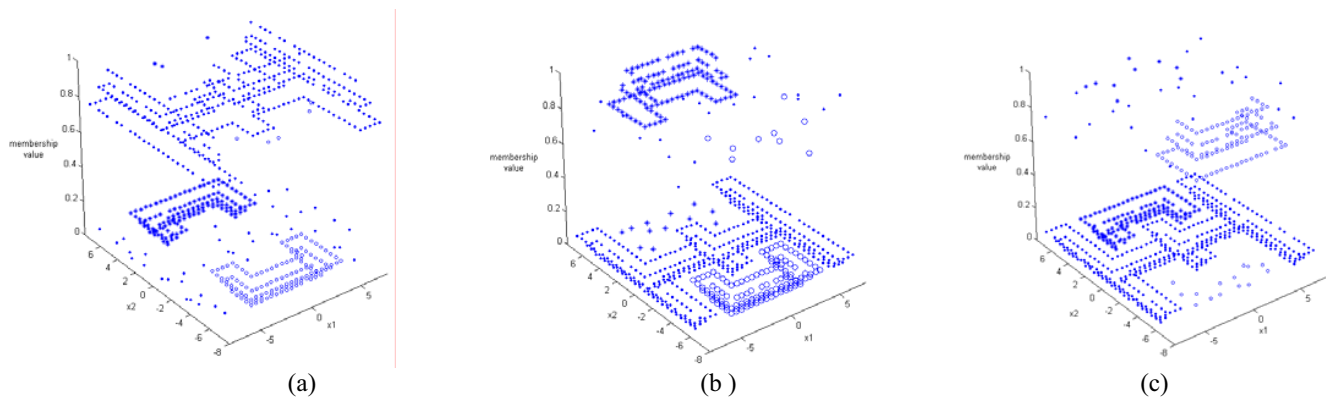


Fig. 1. The distribution of prototypes. (a) class 1; (b) class 2; (c) class3.

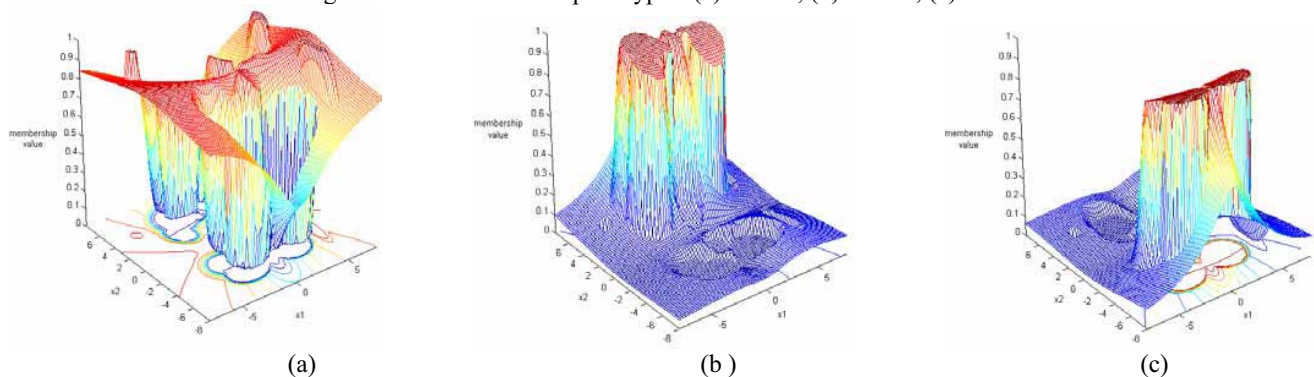


Fig. 2. The decision regions for (a) class 1; (b) class 2; (c) class3.