

Initial Setting of Effective Parameters for Fuzzy Modeling

Hirofumi Miyajima, Noritaka Shigei, and Hiromi Miyajima

Abstract—Many studies on fuzzy modeling(learning of fuzzy inference systems) with vector quantization(VQ) and steepest descend method (SDM) have been made. It is known that these methods are superior in the number of rules(parameters) compared with other learning methods. Most of conventional learning methods using VQ are ones that determine initial assignment of center parameters for membership functions in antecedent part using only input part of learning data, initial assignment of center parameters for membership functions in antecedent part using all learning data and the initial assignment of all parameters in systems using VQ and the generalized inverse matrix(GIM). These methods are ones that determine the initial assignment of parameters in learning process, and any learning data in learning steps of SDM is selected randomly. On the other hand, it is known that many fuzzy rules are needed at or near the places where output changes rapidly in learning data. Therefore, the rate of output change for learning data must be considered. In this paper, we propose learning methods that any data in learning steps of SDM are selected using the probability based on the rate of output change for learning data. Further, the method using GIM to determine the initial assignment of the consequent part of fuzzy rule is also proposed. In order to demonstrate the effectiveness of the proposed methods, numerical simulations for function approximation and pattern classification problems are performed.

Index Terms—Fuzzy Inference Systems, Vector Quantization, Neural Gas Network, Steepest Descent Method, Generalized Inverse Matrix.

I. INTRODUCTION

MANY studies on fuzzy modeling(learning of fuzzy inference systems) have been made [1]–[4]. Their aim is to construct automatically fuzzy inference systems from learning data. Although most of conventional methods are based on steepest descend method(SDM), the obvious drawbacks of them are its large time complexity and getting stuck in a shallow local minimum. Further, there is problems of difficulty dealing with high dimensional spaces [5], [6]. In order to overcome them, some novel methods have been developed, which 1) create fuzzy rules one by one starting from any number of rules, or delete fuzzy rules one by one starting from a sufficiently large number of rules [7], [8], 2) use GA (Genetic Algorithm) and PSO (Particle Swarm Optimization) to determine fuzzy systems [10], 3) use fuzzy inference systems composed of small number of input rule modules, such as SIRMs (Single Input Rule Modules) and DIRMs (Double Input Rule Modules) methods [10], [11], and 4) use a self-organization or a vector quantization technique to

determine the initial assignment of learning parameters [8], [12]. Further, it has been known that a learning method of radial basis function(RBF) networks using VQ and GIM is effective. Specifically, learning methods using VQ and SDM are superior in the number of rules(parameters) compared with other learning methods [2], [12]. Most of conventional learning methods using VQ are ones that determine initial assignment of parameters for membership functions in antecedent part using only input part of learning data. Therefore, we proposed some learning methods to determine initial assignment of center parameters for membership functions in antecedent part using all learning data. Further, we proposed learning methods determining the initial assignment of weight parameters in consequent part [14]–[16]. These methods are ones that determine the initial assignment of learning parameters in learning process, and any learning data is selected randomly in learning steps of SDM [15]–[18]. On the other hand, it is known that many rules are needed at or near the places where output changes rapidly in learning data. Therefore, the rate of change for output data must be considered. Little learning methods selected any data using the probability based on the rate of output change for learning data have been proposed. In this paper, we propose learning methods that any data are selected using the probability based on the rate of output change for learning data in learning process of SDM. Further, the method using GIM to determine initial weights is also proposed. In order to demonstrate the effectiveness of the proposed method, numerical simulations for function approximation and pattern classification problems are performed. In Section II, the conventional learning, VQ and GIM methods are introduced. Further, the probability based on the rate of output change for learning data is introduced. In Section III, new learning methods are proposed. In Section IV, numerical simulations for function approximation and classification problems are performed to show the performance of proposed methods.

II. PRELIMINARIES

A. The conventional fuzzy inference model

The conventional fuzzy inference model using SDM is described [1], [2]. Let $Z_j = \{1, \dots, j\}$ and $Z_j^* = \{0, 1, \dots, j\}$ for the positive integer j . Let R be the set of real numbers. Let $\mathbf{x} = (x_1, \dots, x_m)$ and y^r be input and output data, respectively, where $x_i \in R$ for $i \in Z_m$ and $y^r \in R$. Then the rule of simplified fuzzy inference model is expressed as

$$R_j : \text{if } x_1 \text{ is } M_{1j} \text{ and } \dots \text{ and } x_m \text{ is } M_{mj} \text{ then } y \text{ is } w_j, \quad (1)$$

where $j \in Z_n$ is a rule number, $i \in Z_m$ is a variable number, M_{ij} is a membership function of the antecedent part, and w_j is the weight of the consequent part.

Hirofumi Miyajima is with the Graduate School of Biomedical Sciences, Nagasaki University, Sakamoto, Nagasaki, Japan e-mail: k3768085@kadai.jp.

Noritaka Shigei is with Kagoshima University, e-mail: shigei@eee.kagoshima-u.ac.jp

Hiromi Miyajima is with Kagoshima University, e-mail: miya@eee.kagoshima-u.ac.jp.

A membership value of the antecedent part μ_j for input \mathbf{x} is expressed as

$$\mu_i = \prod_{j=1}^m M_{ij}(x_j). \quad (2)$$

If Gaussian membership function is used, then M_{ij} is expressed as follow:

$$M_{ij}(x_j) = \exp\left(-\frac{1}{2} \left(\frac{x_j - c_{ij}}{b_{ij}}\right)^2\right). \quad (3)$$

, where c_{ij} and b_{ij} are the center and the width values of M_{ij} , respectively. The inference output y^* is calculated by Eq.(4).

$$y^* = \frac{\sum_{i=1}^n \mu_i \cdot w_i}{\sum_{i=1}^n \mu_i}. \quad (4)$$

In order to construct the effective model, the conventional learning is introduced. The objective function E is determined to evaluate the inference error between the desirable output y^r and the inference output y^* .

In this section, we describe the conventional learning algorithm [2].

Let $D = \{(x_1^p, \dots, x_m^p, y^p) | p \in Z_P\}$ and $D^* = \{(x_1^p, \dots, x_m^p) | p \in Z_P\}$ be the set of learning data and the set of input data of D , respectively. The objective of learning is to minimize the following mean square error(MSE):

$$E = \frac{1}{P} \sum_{p=1}^P (y_p^* - y_p^r)^2. \quad (5)$$

, where y_p^* is the inference output for the data x^p .

In order to minimize the objective function E , each parameter $\alpha \in \{c_{ij}, b_{ij}, w_j\}$ is updated based on SDM as follows [1], [2]:

$$\alpha(t+1) = \alpha(t) - K_\alpha \frac{\partial E}{\partial \alpha} \quad (6)$$

where t is iteration time and K_α is a constant. When the Gaussian membership function is used as the membership function, the following relation holds.

$$\frac{\partial E}{\partial c_{ij}} = \frac{\mu_j}{\sum_{j=1}^n \mu_j} \cdot (y^* - y^r) \cdot (w_j - y^*) \cdot \frac{x_j - c_{ij}}{b_{ij}^2} \quad (7)$$

$$\frac{\partial E}{\partial b_{ij}} = \frac{\mu_j}{\sum_{j=1}^n \mu_j} \cdot (y^* - y^r) \cdot (w_j - y^*) \cdot \frac{(x_j - c_{ij})^2}{b_{ij}^3} \quad (8)$$

$$\frac{\partial E}{\partial w_j} = \frac{\mu_j}{\sum_{j=1}^n \mu_j} \cdot (y^* - y^r) \quad (9)$$

The conventional learning algorithm is shown as Fig.1 [1], [2], where n_0 , θ and T_{max} are the initial number of rules, threshold and the maximum number of learning, respectively. Remark that the method is generative one. The method is called learning algorithm A.

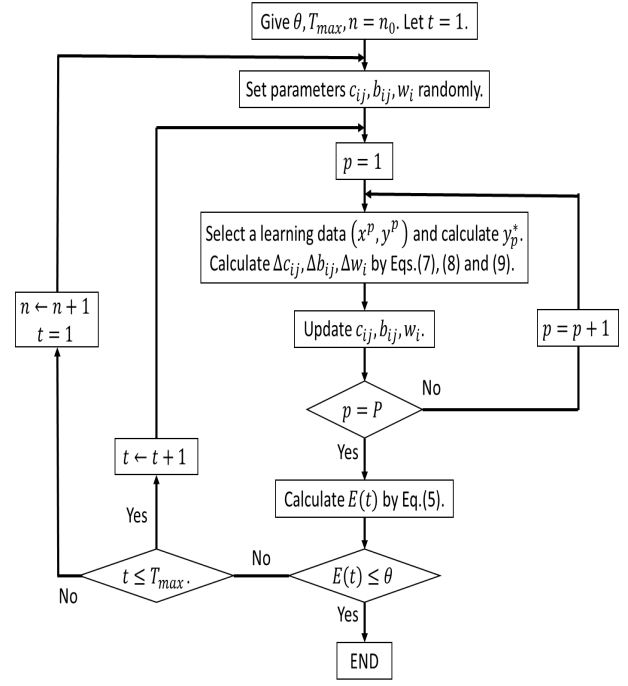


Fig. 1. The flowchart of the conventional learning algorithm

B. Neural gas and K-means methods

Vector quantization techniques encode a data space, e.g., a subspace $V \subseteq R^m$, utilizing only a finite set $C = \{c_i | i \in Z_r\}$ of reference vectors (also called cluster centers), where m and r are positive integers.

Let the winner vector $c_{i(v)}$ be defined for any vector $v \in V$ as follows:

$$i(v) = \arg \min_{i \in Z_r} \|v - c_i\| \quad (10)$$

, where $\|a - b\|$ means the distance between vectors a and b .

From the finite set C , V is partitioned as follows:

$$V_i = \{v \in V | \|v - c_i\| \leq \|v - c_j\| \text{ for } j \in Z_r\} \quad (11)$$

The evaluation function for the partition is defined as follows:

$$E = \sum_{i=1}^r \sum_{v \in V_i} \|v - c_{i(v)}\|^2 \quad (12)$$

For neural gas method [19], the following method is used:

Given an input data vector v , we determine the neighborhood-ranking $(c_{i_0}, \dots, c_{i_{r-1}})$, being the reference vector for which there are k vectors c_j with

$$\|v - c_j\| < \|v - c_{i_k}\| \quad (13)$$

If we denote the number k associated with each vector c_i by $k_i(v, c_i)$, then the adaption step for adjusting the c_i 's is given by

$$\Delta c_i = \varepsilon \cdot h_\lambda(k_i(v, c_i)) \cdot (v - c_i) \quad (14)$$

$$h_\lambda(k_i(v, c_i)) = \exp(-k_i(v, c_i)/\lambda) \quad (15)$$

where $\varepsilon \in [0, 1]$ and $\lambda > 0$. The number λ is called decay constant.

If $\lambda \rightarrow 0$, Eq.(14) becomes equivalent to the K-means method [19]. Otherwise, not only the winner c_{i_0} but the second, third nearest reference vector c_{i_1} , c_{i_2} , etc., are also updated.

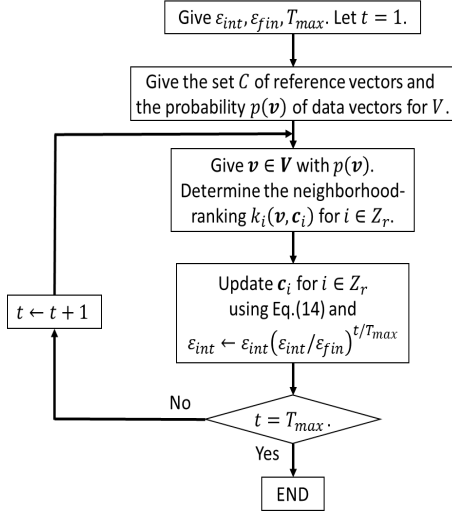


Fig. 2. Neural Gas method

Let $p(v)$ be the probability distribution of data vectors for V . The flowchart of the conventional neural gas algorithm is shown as Fig.2 [19], where ε_{int} , ε_{fin} , θ and T_{max} are learning constants, threshold and the maximum number of learning, respectively. The method is called learning algorithm NG.

If the data distribution $p(v)$ is not given in advance, a stochastic sequence of input data $v(1), v(2), \dots$ which is based on $p(v)$ is given [19].

By using Learning Algorithm NG, learning method of fuzzy systems is shown as follows [12], [13] : In this case, assume that the distribution of learning data is discrete uniform one. Let n_0 be the initial number of rules.

Learning Algorithm B

Step B1 : For learning data D^* , Learning Algorithm NG is performed by using D^* as the set V . The set D^* is encoded by the set C of reference vectors, where $|C| = n_0$.

Step B2 : The set of center parameters of fuzzy rules is set to the set C .

Let

$$b_{ij} = \frac{1}{m_i} \sum_{\mathbf{x}_k \in C_i} (c_{ij} - x_{kj})^2, \quad (16)$$

where C_i and m_i are the i -th cluster for C and the number of learning data for $i \in Z_{n_0}$. Each initial weight w_i is selected randomly.

Step B3 : Learning algorithm A for initial parameters c_{ij} , b_{ij} and w_i are performed.

C. Determination of weights using the generalized inverse method

Let us explain fuzzy inference systems and interpolation problem using the generalized inverse method [2]. This problem can be stated mathematically as follows:

Given the set of learning data $D = \{(\mathbf{x}^p, y_p^r) | p \in Z_P\}$, find a function $f : R^m \rightarrow R$ such that the following conditions are satisfied :

$$f(\mathbf{x}^p) = y_p^r \quad (p \in Z_P) \quad (17)$$

In the case of fuzzy inference system, this problem is solved as follows :

$$y_p = f(\mathbf{x}^p) = \sum_{i=1}^n w_i \phi_{pi}(\|\mathbf{x}^p - \mathbf{c}_i\|) \quad (18)$$

$$\phi_{pi}(\|\mathbf{x}^p - \mathbf{c}_i\|) = \frac{\mu_i(\mathbf{x}^p)}{\sum_{l=1}^n \mu_l(\mathbf{x}^p)} \quad (19)$$

that is,

$$\Phi \mathbf{w} = \mathbf{y} \quad (20)$$

, where $\Phi = (\phi_{pi})$ for $p \in Z_P$ and $i \in Z_n$ is a matrix and $\mathbf{y} = (y_1, \dots, y_r)^T$ and $\mathbf{w} = (w_1, \dots, w_n)^T$ are output and weight vectors, respectively, and T means transportation.

Let $P = n$ and $\mathbf{x}^i = \mathbf{c}_i$. The width parameters are determined by Eq.(16). Then, if $\phi(\cdot)$ is selected as Gaussian function, then the solution of weights \mathbf{w} is obtained as

$$\mathbf{w} = \Phi^{-1} \mathbf{y} \quad (21)$$

Let us consider the case $n < P$. This is the case where the number of rules is sufficiently small for the number of learning data. The optimal solution \mathbf{w}^* that minimizes $E = \|\mathbf{y}^r - \Phi \mathbf{w}\|^2$ can be obtained as follows :

$$\mathbf{w}^+ = \Phi^T \mathbf{y} \text{ and } E_{min} = \|(I - \Psi) \mathbf{y}\|^2 \quad (22)$$

, where $\Phi^+ \triangleq [\Phi^T \Phi]^{-1} \Phi^T$, $\Psi \triangleq \Phi \Phi^T$, I is identify matrix of $P \times P$ and $\mathbf{y} = (y_1^r, \dots, y_P^r)$.

The matrix Φ^+ is called the generalized inverse matrix of Φ . The method using Φ^+ to determine the weights is called the generalized inverse method(GIM) [2].

[Example 1]

Let us perform two-variable logical functions as shown in Table I using fuzzy inference systems with following rules:

$$R_i : \text{if } x_1 \text{ is } M_{i1} \text{ and } x_2 \text{ is } M_{i2} \text{ then } y \text{ is } w_i \quad (23)$$

, where $i \in Z_4$.

Then, the fuzzy inference system with four fuzzy rules is assumed to be of the form:

$$y = \sum_{i=1}^4 \frac{\mu_i w_i}{\sum_{l=1}^4 \mu_l} \quad (24)$$

Let four binary input vectors and the corresponding output vector be $\mathbf{x}_1 = (0, 0)^T$, $\mathbf{x}_2 = (0, 1)^T$, $\mathbf{x}_3 = (1, 0)^T$, $\mathbf{x}_4 = (1, 1)^T$ and $\mathbf{y} = (e_1, e_2, e_3, e_4)^T$ for $e_i \in \{0, 1\}$, respectively(See Table I). Let the center and width parameters of membership functions be $\mathbf{c}_1 = \mathbf{x}_1$, $\mathbf{c}_2 = \mathbf{x}_2$, $\mathbf{c}_3 = \mathbf{x}_3$, $\mathbf{c}_4 = \mathbf{x}_4$ and $b_{ij} = 0.5$ for $i \in Z_4$ and $j \in Z_2$, respectively. Then inference output are obtained as Table I.

In this case, the interpolation matrix Φ may be calculated as

$$\Phi = \begin{bmatrix} 0.778 & 0.105 & 0.105 & 0.012 \\ 0.105 & 0.778 & 0.012 & 0.105 \\ 0.105 & 0.012 & 0.778 & 0.105 \\ 0.012 & 0.105 & 0.105 & 0.778 \end{bmatrix} \quad (25)$$

As a result, weight vectors $\mathbf{w} = (w_1, w_2, w_3, w_4)^T = (0.03, -0.18, -0.18, 1.33)^T$, $(-0.33, 1.18, 1.18, 0.97)^T$ and $(-0.36, 1.36, 1.36, -0.36)^T$ for AND, OR and EXOR operations are obtained, respectively.

TABLE I
FOUR LOGICAL FUNCTIONS

$x \backslash y$	Target output			Inference output		
	AND	OR	EXOR	AND	OR	EXOR
0 0	0	0	0	0.0015	0.0027	0.0012
0 1	0	1	1	0	0.9994	0.9988
1 0	0	1	1	0	0.9994	0.9988
1 1	1	1	0	0.9973	0.9985	0.0012

D. The probability based on the rate of output change for learning data

Learning Algorithm B is a method that determines the initial assignment of fuzzy rules by vector quantization using the set D^* of input for learning data. In this case, the set of output in learning data D is not used to determine the initial assignment of fuzzy rules. In the previous paper, we proposed a method considering both input and output data to determine the initial assignment of fuzzy rules [8].

Based on the improved method of Ref. [8], the probability distribution for D^* is defined as follows : Let D and D^* be the sets of learning data defined in 2.1.

Calculation of the probability for learning data

Step 1 : Give an input data $\mathbf{x}^i \in D^*$, we determine the neighborhood-ranking $(\mathbf{x}^{i_0}, \mathbf{x}^{i_1}, \dots, \mathbf{x}^{i_k}, \dots, \mathbf{x}^{i_{P-1}})$ of the vector \mathbf{x}^i with $\mathbf{x}^{i_0} = \mathbf{x}^i$, \mathbf{x}^{i_1} being closest to \mathbf{x}^i and \mathbf{x}^{i_k} ($k = 0, \dots, P-1$) being the vector \mathbf{x}^i for which there are k vectors \mathbf{x}^j with $\|\mathbf{x}^i - \mathbf{x}^j\| < \|\mathbf{x}^i - \mathbf{x}^{i_k}\|$.

Step 2 : Determine $H(\mathbf{x}^i)$ which shows the rate of change of inclination of the output around output data to input data \mathbf{x}^i , by the following equation:

$$H(\mathbf{x}^i) = \sum_{l=1}^M \frac{|y^i - y^{i_l}|}{\|\mathbf{x}^i - \mathbf{x}^{i_l}\|} \quad (26)$$

, where \mathbf{x}^{i_l} for $l \in Z_M$ means the l -th neighborhood-ranking of \mathbf{x}^i , $i \in Z_P$ and y^i and y^{i_l} are output for input \mathbf{x}^i and \mathbf{x}^{i_l} , respectively. The number M means the range of ranking considering $H(\mathbf{x})$.

Step 3 : Determine the probability $p_M(\mathbf{x}^i)$ for $\mathbf{x}^i \in D^*$ by normalizing $H(\mathbf{x}^i)$.

$$p_M(\mathbf{x}^i) = \frac{H(\mathbf{x}^i)}{\sum_{j=1}^P H(\mathbf{x}^j)} \quad (27)$$

and $\sum_{i=1}^P p_M(\mathbf{x}^i) = 1$.

[Example 2]

Let us explain how to compute $p_M(\mathbf{x})$ using $y = \sin(\pi x_1^3)x_2$ as shown in Fig.4, where $x_1, x_2, y \in [0, 1]$ (See Fig.4).

Assume that four learning data are given as follows :

\mathbf{x}	y
$\mathbf{x}^1 = (0.2, 0.2)$	0.005
$\mathbf{x}^2 = (0.2, 0.8)$	0.020
$\mathbf{x}^3 = (0.8, 0.2)$	0.200
$\mathbf{x}^4 = (0.8, 0.8)$	0.799

The change of the function is rapidly at the point \mathbf{x}^4 and flat at the point \mathbf{x}^1 . Therefore, it is desired to select the point \mathbf{x}^4 with high probability so that many rules assign near this point. Let $M = 2$.

Then, $H(\mathbf{x}^1)$ is calculated as follows :

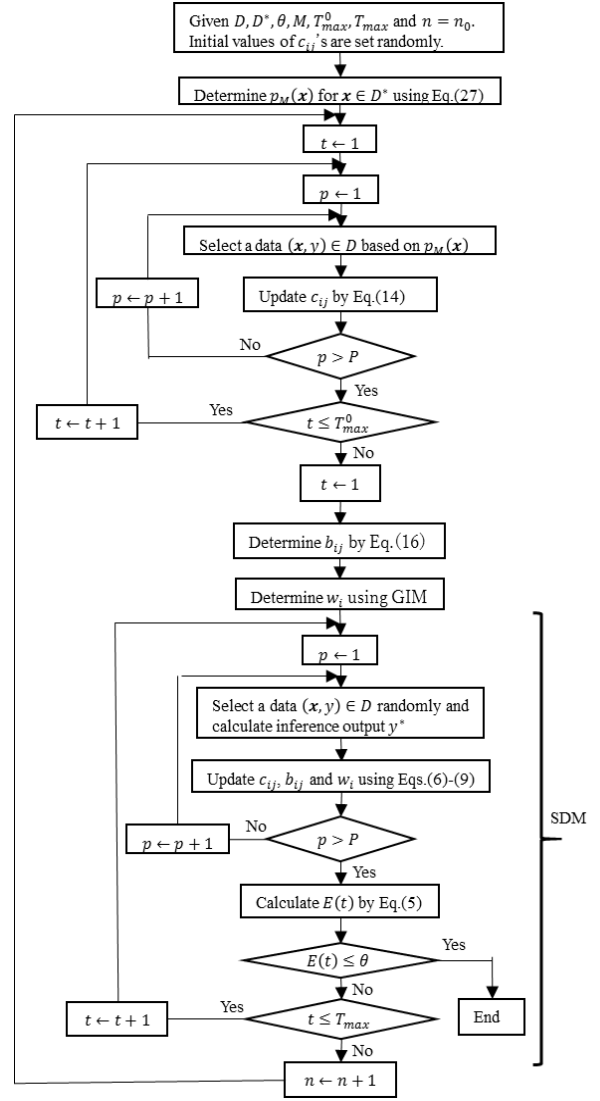


Fig. 3. The flowchart of learning algorithm D

$$H(\mathbf{x}^1) = \frac{|y^1 - y^2|}{\|\mathbf{x}^1 - \mathbf{x}^2\|} + \frac{|y^1 - y^3|}{\|\mathbf{x}^1 - \mathbf{x}^3\|} \quad (28)$$

$$= 0.35$$

because the first and second closest vectors for \mathbf{x}^1 are \mathbf{x}^2 and \mathbf{x}^3 . Likewise, we obtained $H(\mathbf{x}^2) = 1.325$, $H(\mathbf{x}^3) = 1.325$ and $H(\mathbf{x}^4) = 2.3$. From the Eq.(27), each of $p_2(\mathbf{x})$'s is calculated as $p_2(\mathbf{x}^1) = 0.066$, $p_2(\mathbf{x}^2) = 0.25$, $p_2(\mathbf{x}^3) = 0.25$, $p_2(\mathbf{x}^4) = 0.434$.

If $M = 3$, then the following results are obtained : $p_3(\mathbf{x}^1) = 0.143$, $p_3(\mathbf{x}^2) = 0.21$, $p_3(\mathbf{x}^3) = 0.169$, $p_3(\mathbf{x}^4) = 0.477$.

In both cases, the rate of output change for $\mathbf{x}^4 = (0.8, 0.8)$ is large compared to other points, so $p_M(\mathbf{x}^4)$ is large. On the other hand, the probability of each point is 0.25 if data are randomly selected in the case of the conventional method.

The flowchart of learning algorithm D using $p_M(\mathbf{x})$ and GIM is shown in Fig.3.

Then, let us show how the assignment of fuzzy rules changes after learning using Example 2. [Example 3]

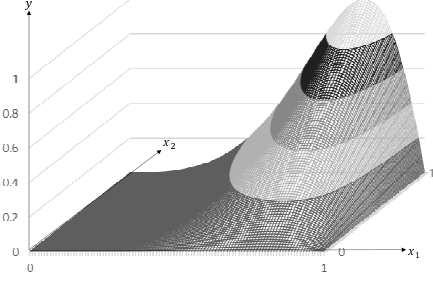


Fig. 4. The figure of $y = \sin(\pi x_1^3)x_2$

Let us consider two cases with $M = 1$ and 100 using $y = \sin(\pi x_1^3)x_2$, where Fig.5(a) and (b) are the initial assignment of center parameters for four fuzzy rules using $p_M(x)$ with $M = 1$ and 100 for the number of learning data $P = 500$, respectively. Fig.5(c) and (d) are the assignment of fuzzy rules after learning for $M = 1$ and 100, respectively. The result shows that the case of $M = 100$ is superior in the assignment of fuzzy rules after learning to the case of $M = 1$.

It is known that many rules are needed at or near palaces where output changes quickly for learning data. The probability $p_M(x)$ is a technique to find the optimum places for the assignment of fuzzy rules. The algorithm D is a heuristic method to find the optimum number of M .

III. THE PROPOSED METHOD

It is shown that learning algorithms A, B and C using VQ and SDM is effective in accuracy and the number of rules to other methods. Further, learning algorithm D using VQ, GIM and SDM is also shown in Ref. [15]. Most of conventional learning methods using VQ and GIM are ones that determine only the initial assignment of parameters in antecedent and consequent parts for membership functions. However, little learning methods using VQ in learning process of SDM have been proposed.

The method using VQ in SDM means that each learning data in SDM is not selected randomly, but selected based on $p_M(x)$. Therefore, each data existing the place where output rapidly changes is more likely to be selected. Let explain it using Example 3.

[Example 4]

Let $M = 2$ in Example 3. Then $p_2(x^1) = 0.066$, $p_2(x^2) = 0.25$, $p_2(x^3) = 0.25$ and $p_2(x^4) = 0.434$, so learning data (x^1, y^1) , (x^2, y^2) , (x^3, y^3) and $(x^4, y^4) \in D$ are selected with the probability 0.066, 0.25, 0.25 and 0.434, respectively.

In this case, output change for (x^4, y^4) and (x^1, y^1) are rapidly and flat as shown in Fig.4, respectively. It seems to assign a lot of fuzzy rules at or near the place (x^4, y^4) .

In this section, we propose four algorithms using $p_M(x)$ in learning steps of SDM corresponding to algorithms A, B, C and D. They are called algorithms(methods) A', B' C' and D'. The algorithm D' is only shown as follows :

Learning Algorithm D'

Step 1 : θ , T_{max}^0 , T_{max} and M are set. Initial values of c_{ij} and b_{ij} are set randomly. $n \leftarrow n_0$.

Step 2 : Let $t = 1$.

Step 3 : Select a data (x^p, y^p) based on $p_M(x^p)$ for $p \in Z_P$.

Step 4 : Update c_{ij} by Eq.(14) : Determination of the center parameters based on $p_M(x)$.

Step 5 : If $t < T_{max}^0$, go to Step 3 with $t \leftarrow t + 1$, otherwise go to Step 6 with $t \leftarrow 1$.

Step 6 : Determine b_{ij} by Eq.(16).

Step 7 : Determine w_i by Eq.(22) : Determination of the initial assignment of the weight w_i by GIM.

Step 8 : Let $p = 1$.

Step 9 : Given a data $(x^p, y_p^r) \in D$ based on $p_M(x)$ for $x \in D$.

Step 10 : Calculate μ_i and y^* by Eqs.(2) and (4).

Step 11 : Update parameters c_{ij} , b_{ij} and w_{ij} by Eqs.(7), (8) and (9).

Step 12 : If $p < P$ then go to Step 8 with $p \leftarrow p + 1$.

Step 13 : If $E < \theta$ or $t > T_{max}$ go to Step 13, otherwise go to Step 8 with $t \leftarrow t + 1$, where E is computed as Eq.(5).

Step 14 : If $E < \theta$, the algorithm terminates, otherwise go to Step 2 with $n \leftarrow n + 1$ and c_{ij} and b_{ij} are set randomly.

The feature of the proposed method is that learning data are selected from the probability $p_M(x)$ in both determining of the initial assignment of parameters and learning steps of SDM(See steps 3 and 9 in learning algorithm D').

In order to compare proposed methods with conventional ones, the following methods are introduced(See Fig.6):

(A) Method A is one based on the conventional algorithm of Fig.1 [1], [2]. Initial parameters of the center c , the width b for the antecedent part and the weight w for the consequent part are set randomly and all parameters are updated by SDM using learning data selected randomly until the inference error becomes sufficiently small.

(B) Method B is known as a learning method of RBF networks [2], [13]. Initial values of c are determined using the set $D^* = \{x^p | p \in Z_P\}$ by VQ and b is computed using c . Initial values of w are randomly selected. Further, all parameters are updated by SDM until the inference error becomes sufficiently small.

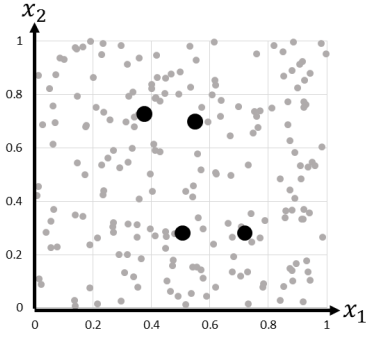
(C) Method C was introduced in Refs. [8], [18]. Initial values of c are determined using the set $D = \{(x^p, y^p) | p \in Z_P\}$ by VQ and b is computed using c . Weight parameters are randomly selected. Further, all parameters are updated by SDM until the inference error becomes sufficiently small.

(D) Method D was introduced in the chapter III. Initial values of c are determined using the set D by VQ and b is computed using c . Initial parameters of w are determined using GIM. Further, all parameters are updated by SDM until the inference error becomes sufficiently small.

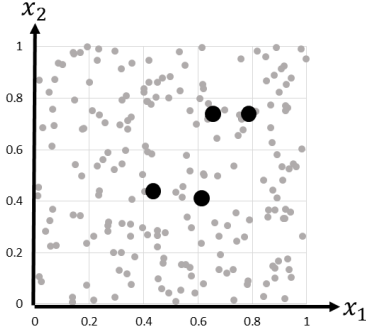
(A') Method A' is the proposed one. Initial parameters of c , b and w are set randomly and all parameters are updated using SDM using learning data based on $p_M(x)$ until the inference error becomes sufficiently small.

(B') Method B' is the proposed one. Initial values of c are determined using the set D^* by VQ and b is computed using c . Weight parameters w are randomly selected. Further, all parameters are updated by SDM based on $p_M(x)$ until the inference error becomes sufficiently small.

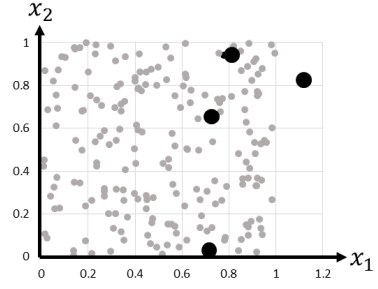
(C') Method C' is the proposed one. Initial values of c are determined using the set D by VQ and b is computed using c . Weight parameters are randomly selected. Further, all parameters are updated by SDM based on $p_M(x)$ until



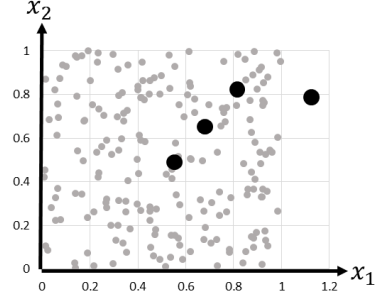
(a)The initial assignment of rules for $M = 1$



(b)The initial assignment of rules for $M = 100$



(c)The assignment of rules after learning for $M = 1$



(d)The assignment of rules after learning for $M = 100$

Fig. 5. The figures (a) and (b) show the initial assignment for $M = 1$ and 100, respectively, where \circ and \bullet mean the places of learning data and center parameters of fuzzy rules. The figures (c) and (d) show the assignment after learning $M = 1$ and 100, respectively.

the inference error becomes sufficiently small.

(D') Method D is the proposed one. Initial values of c are determined using the set D by VQ and b is computed using c . Initial parameters of w are determined using GIM. Further, all parameters are updated by SDM based on $p_M(x)$ until the inference error becomes sufficiently small.

IV. NUMERICAL SIMULATIONS

In order to show the effectiveness of proposed algorithms, simulations of function approximation and classification problems are performed.

A. Function approximation

The systems are identified by fuzzy inference systems. This simulation uses four systems specified by the following functions with 4-dimensional input space $[0, 1]^4$ (Eqs.(29) and (30)) and $[-1, 1]^4$ ((31) and (32)), and one output with the range $[0, 1]$;

$$y = \frac{(2x_1 + 4x_2^2 + 0.1)^2}{37.21} \times \frac{(4 \sin(\pi x_3) + 2 \cos(\pi x_4) + 6)}{12} \quad (29)$$

$$y = \frac{(\sin(2\pi x_1) \times \cos(x_2) \times \sin(\pi x_3) \times x_4 + 1.0)}{2.0} \quad (30)$$

$$y = \frac{(2x_1 + 4x_2^2 + 0.1)^2}{74.42} + \frac{(3e^{3x_3} + 2e^{-4x_4})^{-0.5} - 0.077}{4.68} \quad (31)$$

$$y = \frac{(2x_1 + 4x_2^2 + 0.1)^2}{74.42} + \frac{(4 \sin(\pi x_3) + 2 \cos(\pi x_4) + 6)}{446.52} \quad (32)$$

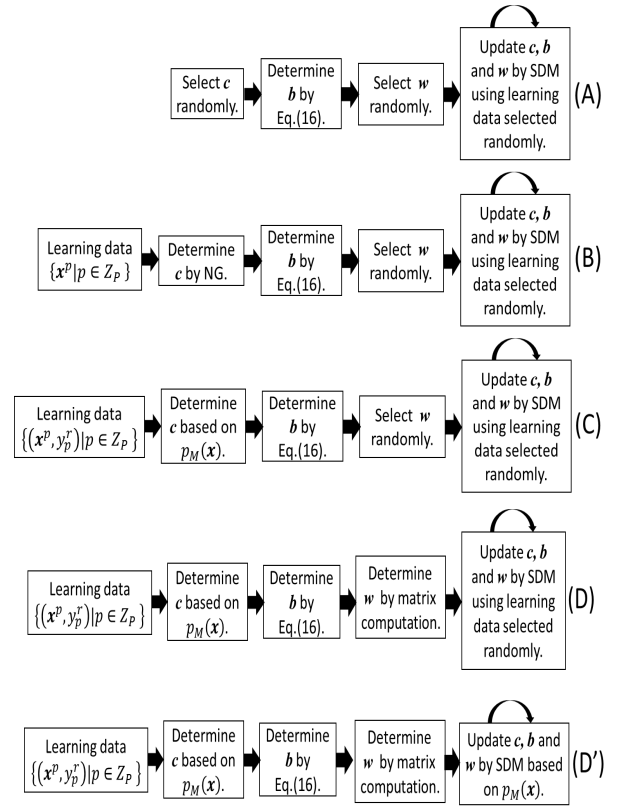


Fig. 6. Concept of conventional and proposed methods, where SDM and NG mean Steepest Descent Method and Neural Gas method. The algorithm D' is only shown, and algorithms A', B' and C' are also defined in the same way.

TABLE II
THE RESULTS FOR FUNCTION APPROXIMATION

		Eq(29)	Eq(30)	Eq(31)	Eq(32)
A	the number of rules	4.2	13.6	7.2	5.1
	MSE for Learning($\times 10^{-4}$)	0.40	0.71	0.43	0.28
	MSE of Test($\times 10^{-4}$)	0.52	1.09	1.00	0.49
B	the number of rules	5.6	14.9	5.2	3.7
	MSE of Learning($\times 10^{-4}$)	0.18	0.77	0.49	0.33
	MSE of Test($\times 10^{-4}$)	0.27	1.42	1.11	0.48
C	the number of rules	4.8	15.6	5.5	4.0
	MSE of Learning($\times 10^{-4}$)	0.21	0.72	0.54	0.88
	MSE of Test($\times 10^{-4}$)	0.34	1.33	0.69	0.53
D	the number of rules	3.1	6.4	4.7	3.2
	MSE of Learning($\times 10^{-4}$)	0.21	0.69	0.56	0.21
	MSE of Test($\times 10^{-4}$)	0.25	0.95	0.68	0.23
A'	the number of rules	3.3	7.8	5.5	3.7
	MSE for Learning($\times 10^{-4}$)	0.22	0.63	0.43	0.27
	MSE of Test($\times 10^{-4}$)	0.30	1.10	0.58	0.33
B'	the number of rules	3.1	7.0	4.9	3.4
	MSE of Learning($\times 10^{-4}$)	0.21	0.66	0.57	0.25
	MSE of Test($\times 10^{-4}$)	0.26	2.30	0.69	0.27
C'	the number of rules	3.2	7.1	4.7	3.1
	MSE of Learning($\times 10^{-4}$)	0.21	0.65	0.60	0.22
	MSE of Test($\times 10^{-4}$)	0.26	0.91	0.73	0.25
D'	the number of rules	3.0	5.9	5.0	3.1
	MSE of Learning($\times 10^{-4}$)	0.25	0.74	0.62	0.24
	MSE of Test($\times 10^{-4}$)	0.35	0.99	0.76	0.26

TABLE III
THE DATASET FOR PATTERN CLASSIFICATION

	Iris	Wine	BCW
The number of data	150	178	683
The number of input	4	13	9
The number of class	3	3	2

As the initial conditions of simulations, $T_{max} = 50000$, $K_c = 0.01$, $K_b = 0.01$, $K_c = 0.1$, $\varepsilon_{init} = 0.1$, $\varepsilon_{fin} = 0.01$, $\lambda = 0.7$, $\theta = 1.0 \times 10^{-4}$ and $M = 200$ are used. The numbers of learning and test data are 512 and 6400, respectively.

Table II shows the result of simulation, where the number of rules, MSE's for learning and test data are shown. In Table II, the number of rules means one when the threshold $\theta = 1.0 \times 10^{-4}$ of inference error is achieved in learning. The result of simulation is the average value from twenty trials. As a result, proposed methods A', B', C' and D' reduce the number of rules compared to conventional methods. Specifically, algorithm D' is superior in the number of rules to other algorithms.

B. Classification problems for UCI database

Iris, Wine and BCW data from UCI database shown in Table III are used as the second numerical simulation [20]. In this simulation, 5-fold cross-validation is used. As the initial conditions for classification problem, $T_{max} = 50000$, $K_c = 0.001$, $K_b = 0.001$, $K_w = 0.05$, $\varepsilon_{init} = 0.1$, $\varepsilon_{fin} = 0.01$ and $\lambda = 0.7$ are used. Further, $M = 100$ and $\theta = 1.0 \times 10^{-2}$ for Iris and Wine and $M = 200$ and $\theta = 2.0 \times 10^{-2}$ for BCW are used.

Table IV shows the result of classification problem. In Table IV, the number of rules, RM's for learning and test data are shown, where RM means the rate of misclassification. As a result, it is shown that proposed methods A', B', C' and D' reduce the number of rules in classification problem.

TABLE IV
THE RESULT FOR PATTERN CLASSIFICATION

		Iris	Wine	BCW
A	the number of rules	3.4	7.8	14.4
	RM for Learning(%)	3.0	1.4	1.6
	RM of Test(%)	3.3	10.3	4.3
B	the number of rules	2.0	20.8	26.0
	RM of Learning(%)	3.3	13.6	2.2
	RM of Test(%)	3.3	16.6	3.5
C	the number of rules	3.4	7.4	9.6
	RM of Learning(%)	2.8	2.1	2.0
	RM of Test(%)	4.7	5.1	4.6
D	the number of rules	2.0	2.0	2.3
	RM of Learning(%)	2.5	1.5	1.5
	RM of Test(%)	4.0	8.1	4.0
A'	the number of rules	2.2	2.2	3.7
	RM for Learning(%)	2.7	1.5	1.5
	RM of Test(%)	3.6	7.7	4.1
B'	the number of rules	2.2	2.6	2.5
	RM of Learning(%)	2.4	1.3	1.7
	RM of Test(%)	3.9	7.7	4.0
C'	the number of rules	2.2	2.5	2.6
	RM of Learning(%)	2.4	1.4	1.6
	RM of Test(%)	3.8	8.2	4.0
D'	the number of rules	2.0	2.0	2.3
	RM of Learning(%)	2.5	1.6	1.6
	RM of Test(%)	3.7	8.2	4.1

Let us consider the reason why we can get the good result by using the probability $p_M(x)$. In the conventional learning method, parameters are updated by any data selected randomly from the set of learning data. In the proposed method, parameters are updated by any data selected from the probability $p_M(x)$. The function $p_M(x)$ is determined based on output change for the set of learning data, so many fuzzy rules are likely to generate at or near the places where output change is large for the set of learning data. For example, if the number of learning time is 100 and $p_M(x^0) = 0.5$, then learning data x^0 is selected 50 times from the set of learning data in learning. As a result, membership functions are likely to generate at or near the places where output change is large for the set of learning data. The probability $p_M(x)$ is considered as a method to improve the local search of SDM.

However, there is no difference in capability between algorithms D and D'. This means that the effect of $p_M(x)$ on SDM is limited. In other words, after adjusting the center parameters of rules and the weight of the consequent part as initial setting, even if $p_M(x)$ or randomly selection for SDM is used, the capability may not change. This result needs to be verified by other simulations.

V. CONCLUSION

In this paper, we propose learning methods that any data in learning steps of SDM is selected based on a probability based on the rate of output change for learning data. It means that each learning data in SDM is not selected randomly, but selected based on the probability for output change. Therefore, each data existing the place where output rapidly changes is more likely to be selected. In order to demonstrate the effectiveness of the proposed method, numerical simulations for function approximation and pattern classification problems were performed.

In the future work, we will propose more effective learning algorithm using the probability of output change for learning data compared to other methods and consider to apply the proposed method to learning of neural network.

REFERENCES

- [1] B. Kosko, *Neural Networks and Fuzzy Systems, A Dynamical Systems Approach to Machine Intelligence*, Prentice Hall, Englewood Cliffs, NJ, 1992.
- [2] M.M. Gupta, L. Jin and N. Homma, *Static and Dynamic Neural Networks*, IEEE Press, 2003.
- [3] C. Lin and C. Lee, *Neural Fuzzy Systems*, Prentice Hall, PTR, 1996.
- [4] B. Liu, *Theory and Practice of Uncertain Programming, Studies in Fuzziness and Soft Computing*, Vol. 239, Springer, 2009.
- [5] J. Casillas, O. Cordon, F. Herrera and L. Magdalena, *Accuracy Improvements in Linguistic Fuzzy Modeling, Studies in Fuzziness and Soft Computing*, Vol. 129, Springer, 2003.
- [6] S. M. Zhoua and J. Q. Ganb, *Low-level interpretability and high-level interpretability: a unified view of data-driven interpretable fuzzy system modeling*, *Fuzzy Sets and Systems* 159, pp.3091-3131, 2008.
- [7] H. Nomura, I. Hayashi and N. Wakami, *A Learning Method of Simplified Fuzzy Reasoning by Genetic Algorithm*, *Proc. of the Int. Fuzzy Systems and Intelligent Control Conference*, pp.236-245, 1992.
- [8] K. Kishida and H. Miyajima, *A Learning Method of Fuzzy Inference Rules using Vector Quantization*, *Proc. of the Int. Conf. on Artificial Neural Networks*, Vol.2, pp.827-832, 1998.
- [9] O. Cordon, *A historical review of evolutionary learning methods for Mamdani-type fuzzy rule-based systems, Designing interpretable genetic fuzzy systems*, *Journal of Approximate Reasoning*, 52, pp.894-913, 2011.
- [10] N. Yubazaki, J. Yi and K. Hirota, *SIRMS(Single Input Rule Modules) Connected Fuzzy Inference Model*, *J. Advanced Computational Intelligence*, 1, 1, pp.23-30, 1997.
- [11] H. Miyajima, N. Shigei and H. Miyajima, *Fuzzy Inference Systems Composed of Double-Input Rule Modules for Obstacle Avoidance Problems*, *IAENG International Journal of Computer Science*, Vol. 41, Issue 4, pp.222-230, 2014.
- [12] K. Kishida, H. Miyajima, M. Maeda and S. Murashima, *A Self-tuning Method of Fuzzy Modeling using Vector Quantization*, *Proceedings of FUZZ-IEEE'97*, pp397-402, 1997.
- [13] S. Fukumoto, H. Miyajima, N. Shigei and K. Uchikoba, *A Decision Procedure of the Initial Values of Fuzzy Inference System Using Counterpropagation Networks*, *Journal of Signal Processing*, Vol.9, No.4, pp.335-342, 2005.
- [14] H. Miyajima, N. Shigei, K. Kishida, Y. Akiyoshi and H. Miyajima, *An Improved Learning Algorithm of Fuzzy Inference Systems using Vector Quantization*, *Advanced in Fuzzy Sets and Systems*, Vol.20, pp.155-175, 2015.
- [15] H. Miyajima, N. Shigei, H. Miyajima, "Fast Learning Algorithm for Fuzzy Inference Systems Using Vector Quantization," *Lecture Notes in Engineering and Computer Science: Proceedings of The International MultiConference of Engineers and Computer Scientists 2016, IMECS 2016*, 16-18 March, 2016, Hong Kong, pp1-6.
- [16] H. Miyajima, N. Shigei and H. Miyajima, *The Ability of Learning Algorithms for Fuzzy Inference Systems using Vector Quantization*, *ICONIP 2016, LNCS 9950*, pp.479, Kyoto, October, 2016.
- [17] W. Pedrycz, H. Izakian, *Cluster-Centric Fuzzy Modeling*, *IEEE Trans. on Fuzzy Systems*, Vol. 22, Issue 6, pp. 1585-1597, 2014.
- [18] H. Miyajima, N. Shigei and H. Miyajima, "Fuzzy Modeling using Vector Quantization based on Input and Output Learning Data", *Lecture Notes in Engineering and Computer Science: Proceedings of The International MultiConference of Engineers and Computer Scientists 2017*, 15-17 March, 2017, Hong Kong, pp.1-6.
- [19] T. M. Martinetz, S. G. Berkovich and K. J. Schulten, *Neural Gas Network for Vector Quantization and its Application to Time-series Prediction*, *IEEE Trans. Neural Network*, 4, 4, pp.558-569, 1993.
- [20] *UCI Repository of Machine Learning Databases and Domain Theories*, <ftp://ftp.ics.uci.edu/pub/machinelearning-Databases>.