



ELSEVIER

Available online at www.sciencedirect.com

SCIENCE @ DIRECT®

Nonlinear Analysis 63 (2005) e877–e889

**Nonlinear
Analysis**

www.elsevier.com/locate/na

Parallelization of Fuzzy ARTMAP to improve its convergence speed: The network partitioning approach and the data partitioning approach

José Castro^a, Michael Georgiopoulos^{a,*}, Jimmy Secretan^a,
Ronald F. DeMara^a, Georgios Anagnostopoulos^b, Avelino Gonzalez^a

^aDepartment of Electrical and Computer Engineering, University of Central Florida, Orlando,
FL 32816-2786, USA

^bDepartment of Electrical and Computer Engineering, Florida Institute of Technology, FL 32816 - 2786, USA

Abstract

One of the properties of FAM, which can be both an asset and a liability, is its capacity to produce new neurons (templates) on demand to represent classification categories. This property allows FAM to automatically adapt to the database without having to arbitrarily specify network structure. We provide two methods for speeding up the FAM algorithm. The first one, referred to as the *data partitioning* approach, partitions the data into subsets for independent processing. The second one, referred to as the *network partitioning* approach, uses a pipeline to distribute the work between processes during training. We provide experimental results on a Beowulf cluster of workstations for both approaches that confirm the speedup of the modifications.

© 2005 Elsevier Ltd. All rights reserved.

Keywords: Fuzzy ARTMAP; BEOWULF parallel processing; Data partitioning; Network partitioning

1. Introduction

Neural networks have been used extensively and successfully to address a wide variety of problems. As computing capacity and electronic databases grow, there is an increasing need to process considerably larger datasets. In this context, the algorithms of choice tend

* Corresponding author. Tel.: +1 407 823 5338; fax: +1 407 823 5835.

E-mail address: michaelg@mail.ucf.edu (M. Georgiopoulos).

to be ad-hoc algorithms, or tree-based algorithms such as CART [2] and C4.5 [7]. Variations of these tree learning algorithms, such as SPRINT [8] and SLIQ [5] have been successfully adapted to handle very large data sets.

Neural network algorithms have, for some applications, prohibitively high convergence times. Even one of the fastest neural network algorithms, the Fuzzy ARTMAP (FAM) algorithm, tends to lag in convergence time as the size of the network grows. The FAM algorithm corresponds to a family of neural network architectures introduced by Carpenter et al. [3,4], and has proven to be one of the premier neural network architectures for classification problems.

Some of the advantages that FAM has, compared to other neural network classifiers are that it learns the required task fast, it has the capability to do on-line learning, and its learning structure allows the explanation of the answers that the neural network produces.

One of FAM's properties which is a mixed blessing is its capacity to produce new neurons (*templates*) on demand to represent classification categories. This property allows FAM to automatically adapt to the database without having to arbitrarily and a priori specify its network structure, but it also has the undesirable side effect that for large databases it can produce a large network size that can dramatically slow down the algorithm's training time. It would be desirable to have a method capable of keeping FAM's convergence time manageable, without affecting the generalization performance of the network or its resulting size when the training is completed.

In this paper, we propose two partitioning approaches for the FAM algorithm to be used in a parallel setting. The *network partitioning* approach and the *data partitioning* approach. Our research on the proposed data partitioning approach for FAM has shown that it dramatically reduces the training time of FAM, by partitioning the training set, without a significant effect on the classification (generalization) performance of the network. For the network partitioning approach, a pipelined FAM implementation is proposed that speeds up the algorithm linearly with respect to the number of processors used in the pipeline.

This paper is organized as follows: first we review the Fuzzy ARTMAP architecture and functionality, then we examine the computational complexity of FAM and analyze how and why a data partitioning and/or a network partitioning approach can considerably reduce the algorithm's training time. Then, we discuss in detail the data partitioning of FAM and provide a network partitioning parallel pipelined algorithm for FAM. Furthermore, we discuss some results of correctness on the network partitioning pipelined FAM approach. Finally, experimental results are presented that illustrate the merit of our approaches. The data-partitioned FAM and the pipelined FAM were implemented on a Beowulf cluster of workstations. We close the paper with a summary of the findings and suggestions for further research.

2. Fuzzy ARTMAP

2.1. The Fuzzy ARTMAP architecture

The Fuzzy ARTMAP architecture consists of four layers or fields of nodes (see Fig. 1). The layers that are worth describing are the *input layer* F_1^a , the *category representation*

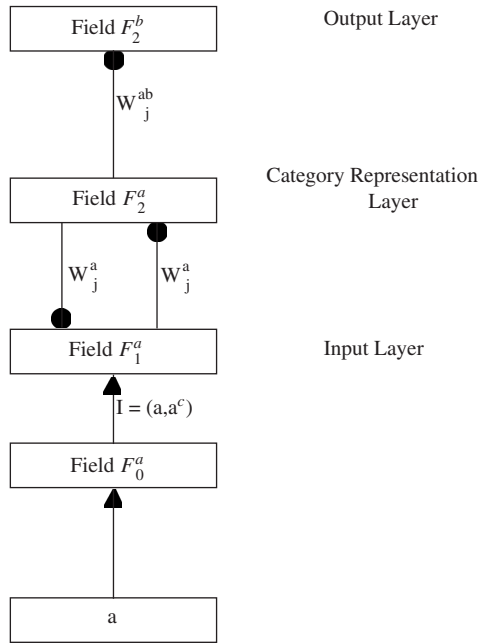


Fig. 1. Fuzzy ARTMAP diagram.

layer F_2^a , and the *output layer* F_2^b . The input layer of Fuzzy ARTMAP is the layer where an input vector of dimensionality $2M_a$ of the following form is applied:

$$\mathbf{I} = (\mathbf{a}, \mathbf{a}^c) = (a_1, a_2, \dots, a_{M_a}, a_1^c, a_2^c, \dots, a_{M_a}^c), \tag{1}$$

$$a_i^c = 1 - a_i, \quad \forall i \in \{1, 2, \dots, M_a\}. \tag{2}$$

The assumption here is that the input vector \mathbf{a} is such that each one of its components lies in the interval $[0, 1]$.

The layer F_2^a of Fuzzy ARTMAP is referred to as the *category representation layer*, because this is where categories (or groups) of input patterns are formed. Finally, the output layer is the layer that produces the outputs of the network. An output of the network represents the output to which the input applied at the input layer of FAM is supposed to be mapped to.

There are two sets of weights worth mentioning in FAM. The first set of weights are weights from F_2^a to F_1^a , designated as w_{ij}^a , $1 \leq j \leq N_a$, $1 \leq i \leq M_a$, and referred to as top-down weights. The vector of weights

$$\mathbf{w}_j^a = (w_{j1}^a, w_{j2}^a, \dots, w_{j,2M_a}^a)$$

is called a *template*. Its functionality is to represent the group of input patterns that chose node j in the category representation layer of Fuzzy ARTMAP as their representative node. The second set of weights worth mentioning, are weights that emanate from every node j

in the category representation layer to every node k in the output layer. These weights are designated as W_{jk}^{ab} (called inter-ART weights). The vector of inter-ART weights emanating from every node j in Fuzzy ARTMAP

$$\mathbf{W}_j^{ab} = (W_{j1}^{ab}, W_{j2}^{ab}, \dots, W_{j,N_b}^{ab})$$

corresponds to the output pattern that this node j is mapped to.

Fuzzy ARTMAP can operate in two distinct phases: the *training phase* and the *performance phase*. The training phase of Fuzzy ARTMAP can be described as follows: Given a list of input/output pairs,

$$\{(\mathbf{I}^1, \mathbf{O}^1), (\mathbf{I}^2, \mathbf{O}^2), \dots, (\mathbf{I}^P, \mathbf{O}^P)\}.$$

We want to train Fuzzy ARTMAP to map every input pattern of the training list to its corresponding output pattern. To achieve the aforementioned goal, we present the training list to Fuzzy ARTMAP architecture repeatedly. That is, we present \mathbf{I}^1 to F_1^a , \mathbf{O}^1 to F_2^b , \mathbf{I}^2 to F_1^a , \mathbf{O}^2 to F_2^b , and finally \mathbf{I}^P to F_1^a , and \mathbf{O}^P to F_2^b . We present the training list to Fuzzy ARTMAP as many times as it is necessary for Fuzzy ARTMAP to correctly classify all these input patterns. The task is considered accomplished (i.e., the learning is complete) when the weights do not change during a list presentation. The aforementioned training scenario is called *off-line* learning. The performance phase of Fuzzy ARTMAP works as follows: Given a list of input patterns, such as $\tilde{\mathbf{I}}^1, \tilde{\mathbf{I}}^2, \dots, \tilde{\mathbf{I}}^S$, we want to find the Fuzzy ARTMAP output produced when each one of the aforementioned test patterns is presented at its F_1^a layer. In order to achieve the aforementioned goal, we present the test list to the trained Fuzzy ARTMAP architecture and we observe the network's output.

The operation of Fuzzy ARTMAP is affected by two network parameters, the choice parameter β_a , and the baseline vigilance parameter $\bar{\rho}_a$. The choice parameter takes values in the interval $(0, \infty)$, while the baseline vigilance parameter assumes values in the interval $[0, 1]$. Both of these parameters affect the number of nodes created in the category representation layer of Fuzzy ARTMAP. Higher values of β_a and $\bar{\rho}_a$ create more nodes in the category representation layer of Fuzzy ARTMAP, and consequently produce less compression of the input patterns. There are two other network parameter values in Fuzzy ARTMAP that are worth mentioning. The vigilance parameter ρ_a and the number of nodes N_a in the category representation layer of Fuzzy ARTMAP. The vigilance parameter ρ_a takes value in the interval $[\bar{\rho}_a, 1]$ and its initial value is set to be equal to $\bar{\rho}_a$. The number of nodes N_a in the category representation layer of Fuzzy ARTMAP increases, while training the network, and it corresponds to the number of committed nodes in Fuzzy ARTMAP plus one uncommitted node.

2.2. The Fuzzy ARTMAP learning algorithm

Prior to initiating the training phase of Fuzzy ARTMAP, the top-down weights (the w_{ij}^a 's) are chosen equal to 1, and the inter-ART weights (the W_{jk}^{ab} 's) are chosen equal to 0. There are three major operations that take place during the presentation of a training input/output pair (e.g., $(\mathbf{I}^r, \mathbf{O}^r)$) to Fuzzy ARTMAP, and they are described below. One of the specific operands involved in all of these operations is the *fuzzy min operand*, designated by the

symbol \wedge . Actually, the fuzzy min operation of two vectors \mathbf{x} , and \mathbf{y} , designated as $\mathbf{x} \wedge \mathbf{y}$, is a vector whose components are equal to the minimum of components of \mathbf{x} and \mathbf{y} . Another specific operand involved in these equations is designated by the symbol $|\bullet|$. In particular, $|\mathbf{x}|$ is the size of a vector \mathbf{x} and is defined to be the sum of its components.

Operation 1. Calculation of bottom-up inputs to every node j in F_2^a as follows:

$$T_j^a = \frac{|\mathbf{I}^r \wedge \mathbf{w}_j^a|}{|\mathbf{w}_j^a| + \beta_a}. \tag{3}$$

After calculation of the bottom up inputs the node j_{\max} with the maximum bottom-up input is chosen.

Operation 2. The node j_{\max} with the maximum bottom up input is examined to determine whether it passes the vigilance criterion. A node passes the vigilance criterion if the following condition is met:

$$\frac{|\mathbf{I}^r \wedge \mathbf{w}_{j_{\max}}^a|}{|\mathbf{I}^r|} \geq \rho_a. \tag{4}$$

If the vigilance criterion is satisfied we proceed with operation 3 otherwise node j_{\max} is disqualified and we find the next node in F_2^a that maximizes the bottom-up input. Eventually we will end up with a node j_{\max} that maximizes the bottom-up input and passes the vigilance criterion.

Operation 3. This operation is implemented only after we have found a node j_{\max} that maximizes the bottom-up input and passes the vigilance criterion. Operation 3 determines whether this node j_{\max} passes the prediction test. The prediction test checks if the inter-ART weight vector emanating from node j_{\max} i.e.

$$\mathbf{W}_{j_{\max}}^{ab} = (W_{j_{\max}1}^{ab}, W_{j_{\max}2}^{ab}, \dots, W_{j_{\max}N_b}^{ab})$$

matches exactly the desired output vector \mathbf{O}^r (if it does, this is referred to as *passing the prediction test*). If the node does *not* pass the prediction test, the vigilance parameter ρ_a is increased to the level of

$$\rho_a \leftarrow \frac{|\mathbf{I}^r \wedge \mathbf{w}_{j_{\max}}^a|}{|\mathbf{I}^r|} + \varepsilon,$$

where ε is a very small number. Node j_{\max} is disqualified, and the next F_2^a node that maximizes the bottom-up input and passes the new vigilance is chosen. If, on the other hand, node j_{\max} passes the predictability test, the weights in Fuzzy ARTMAP are modified as follows:

$$\mathbf{w}_{j_{\max}}^a \leftarrow \mathbf{w}_{j_{\max}}^a \wedge \mathbf{I}^r, \quad \mathbf{W}_{j_{\max}}^{ab} \leftarrow \mathbf{O}^r. \tag{5}$$

Fuzzy ARTMAP training is considered complete if and only if after repeated presentations of all training input/output pairs to the network, where Operations 1–3 are recursively

```

FAM-LEARNING-PHASE(Patterns,  $\bar{\rho}_a, \beta_a, \text{maxEpochs}$ )
1  templates  $\leftarrow \{\}$ 
2  iter  $\leftarrow 0$ 
3  repeat
4      modified  $\leftarrow \text{false}$ 
5      for each I in Patterns
6          do  $\rho_a \leftarrow \bar{\rho}_a$ 
7              LEARN-PATTERN(I, templates,  $\rho_a, \beta_a$ )
8              iterations  $\leftarrow \text{iterations} + 1$ 
9
10     until (iter = maxEpochs) or (modified = false)
11 return templates

```

Where the procedure LEARN-PATTERN is:

```

LEARN-PATTERN(I, templates,  $\rho_a, \beta_a$ )
1  repeat
2      status  $\leftarrow \text{FoundNone}$ 
3      jmax  $\leftarrow \text{GET-MAX-ARG}(\text{I}, \text{templates}, \rho_a, \beta_a)$ 
4      if status = FoundOne
5          then if class(I) = class( $\mathbf{w}_{j_{\text{max}}}^a$ )
6              then
7                  status  $\leftarrow \text{ThisIsIT}$ 
8              else
9                  status  $\leftarrow \text{Matchtracking}$ 
10                  $\rho \leftarrow \rho(\text{I}, \mathbf{w}_{j_{\text{max}}}^a) + \varepsilon$ 
11 until status  $\neq \text{Matchtracking}$ 

```

Fig. 2. Fuzzy ARTMAP algorithm.

applied for every input/output pair, we find ourselves in a situation where a complete cycle through all the input/output pairs produced no weight changes. In some databases, noise in the data may create over-fitting when we repeatedly present the training input/output pairs; so a single pass over the training set may be preferable. A single pass over the training data is also applied when Fuzzy ARTMAP is used in an on-line training environment.

In the performance phase of Fuzzy ARTMAP only Operations 1 and 2 are implemented for every input pattern presented to Fuzzy ARTMAP. By registering the network output to every test input applied to Fuzzy ARTMAP, and by comparing it to the desired output, we can calculate the network's performance (i.e., network's misclassification error).

3. Fuzzy ARTMAP time complexity analysis

3.1. Online Fuzzy ARTMAP time complexity

The pseudo-code for the FAM algorithm is shown in Fig. 2. The FAM algorithm tests every input pattern **I** in the training set against every template \mathbf{w}_j^a in F_2^a at least once. Let us call T the average number of times that the inner searching loop is executed for each input

pattern, and call it the *matchtracking factor*. Then the number of times that a given input pattern \mathbf{I} passes through the code will be

$$Time(\mathbf{I}) = O(\Gamma \times |templates|). \tag{6}$$

If we assume that the number of templates does not change during training it is easy to see that the time complexity of the algorithm is

$$Time(FAM) = O(\Gamma \times P \times |templates|). \tag{7}$$

For a fixed type of database the FAM algorithm achieves a certain *compression ratio*. This means that the number of templates created is actually a fraction of the number of patterns P in the training set. Let us call this compression ratio κ so that

$$|templates| = \kappa P \tag{8}$$

and thus the FAM time complexity is given by the formula

$$O(FAM) = O(\Gamma P \kappa P) = O(\kappa \Gamma P^2). \tag{9}$$

3.2. Off-line Fuzzy ARTMAP time complexity

Without any previous assumptions the complexity of the off-line FAM algorithm would be

$$O \left(\underbrace{\sum_{i=1}^{epochs} \kappa_i \Gamma_i P^2}_A \right), \tag{10}$$

where Γ_i is the matchtracking factor for the i th epoch, and κ_i is the compression ratio of the i th epoch.

We can simplify A , in Eq. (10), by making the reasonable assumptions that: (a) the number of epochs is dependent on the number of training patterns, and (b) Γ_i does not vary from epoch to epoch and (c) $\kappa_i = \kappa^i$. Condition (c) basically means that the amount of templates created decreases geometrically from epoch to epoch. Using these assumptions we can state the formula as:

$$\begin{aligned} A &= \sum_{i=1}^{epochs} \kappa^i \Gamma P^2 = \Gamma P^2 \sum_{i=1}^{epochs} \kappa^i \\ &= \Gamma P^2 \kappa \left(\frac{1 - \kappa^{epochs(P)}}{1 - \kappa} \right). \end{aligned} \tag{11}$$

If we use the stopping criterion where learning stops when $\kappa^{epochs} P = 1$, then we can end up with a computable formula for the term A of Eq. (10). In particular, it can be

shown that

$$A = \kappa \Gamma P^2 \left(\frac{1 - \kappa^{\lceil -\ln(P)/\ln(\kappa) \rceil}}{1 - \kappa} \right), \quad (12)$$

where the number of templates created is given by the following equation:

$$\kappa P \left(\frac{1 - \kappa^{\lceil -\ln(P)/\ln(\kappa) \rceil}}{1 - \kappa} \right). \quad (13)$$

We can see from this formula that the convergence of the FAM algorithm is quadratic with respect to P , the number of patterns in the training set. This relationship suggests that it may be beneficial to use methods that reduce the size of the training set in FAM to speed up the convergence of the algorithm, and that by doing so we may get a quadratic speedup on the time complexity of the algorithm. The following sections are motivated by this partitioning, one reduces the parameter P by partitioning the input set of FAM, the other reduces the number of templates κP by doing network partitioning.

4. Partitioned FAM

4.1. Data partitioned boxing FAM

We chose to follow a data-partitioning approach that allows reducing the number of patterns presented to FAM, and consequently the number of templates created in FAM. Through the data-partitioning scheme that we are proposing the dataset is split into smaller datasets, each one of which trains a different FAM network. This approach has the advantage that it can reduce the computational complexity of the algorithm and lends itself well to parallel implementation.

Under the fairly reasonable assumption that the number of templates in the ART neural network is proportional to the number of training patterns, we can state that the ART convergence time is quadratic (i.e., $O(P^2)$), where as before, P is the number of patterns in the data set. It is therefore obvious that partitioning the data set will result in a significant improvement for the FAM convergence time. For instance, if we divide the data set into p equal partitions of size P/p , the partitioning will theoretically produce a speedup in the order of p^2 . So we should theoretically expect a quadratic improvement in speed as we increase the number of partitions in the data set.

Our approach is inspired by the projective clustering approach proposed and implemented in [6], and by the natural properties of the FAM algorithm. We project the data contained in the M_a -dimensional hypercube to \hat{M}_a dimensions, where $\hat{M}_a \ll M_a$, and we partition the data set by a grid in this projected space. If, for example, we use $\hat{M}_a = 3$ dimensions and a grid size of 10 divisions per dimension, we would divide the 3-dimensional space, on which we projected the original data, into 1000 boxes of side length equal to 0.1. If each set of data within a box trains a different FAM architecture (partitioned FAM) we are guaranteed that the number of templates created by each such FAM is not going to be large

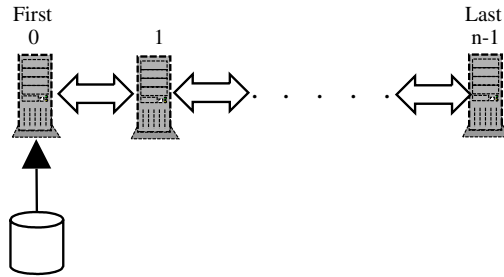


Fig. 3. Pipeline structure.

(especially if the number of boxes is large). It is likely though that the total number of templates created by the partitioned FAM is larger than the number of templates created by the non-partitioned FAM. Furthermore, classification performance may also be negatively affected by this partitioning approach. To avoid this unwelcome side effect, we take advantage of the natural partitioning imposed by the FAM algorithm. We know that templates in FAM are geometrically represented by hyper-rectangles. Furthermore each hyper-rectangle size is restricted by the vigilance parameter $\bar{\rho}_a$ and by the dimensionality M_a of the input patterns. In particular

$$size(\mathbf{w}_j^a) \leq M_a(1 - \bar{\rho}_a). \tag{14}$$

If we assume that templates grow more or less evenly across every dimension, then by choosing boxes in the projected space of size $(1 - \bar{\rho}_a)$, across every dimension, we are actually restricting the template size to the size that the FAM algorithm already enforces. This partitioning approach is most effective when the value of the vigilance parameter is large, and this is the case when the number of templates grows the most, and tends to slow down the training of the algorithm.

4.2. Network partitioned pipelined FAM

In this second approach, we chose to do network partitioning to bring down the complexity of the FAM algorithm. This approach requires the use of a pipeline of processors which exchange input patterns and templates through the pipeline. The pipeline structure is depicted in Fig. 3.

One of the problems presented in pipelining the FAM algorithm is its use of matchtracking: a condition that if met induces the network to retrain a pattern and compare it again to all the templates in the network. This situation poses serious issues to a pipeline design, since when a pattern is subjected to a matchtracking operation we must flush all the elements of the pipeline to guarantee equivalence to the sequential algorithm. Obviously this approach is inefficient, so we chose to implement a variant of FAM, called no-matchtracking FAM, developed by Anagnostopoulos (see [1]). The modified FAM algorithm is depicted in Fig. 4.

```

FAM-NO-MATCHTRACKING-LEARNING( $\{\mathbf{I}^1, \dots, \mathbf{I}^P\}, \rho_a, \beta_a$ )
1  $\mathbf{w}_0 \leftarrow \underbrace{(1, 1, \dots, 1)}_{2M_a}$ 
2  $templates \leftarrow \{\mathbf{w}_0\}$ 
3 for each  $\mathbf{I}^r$  in  $\{\mathbf{I}^1, \mathbf{I}^2, \dots, \mathbf{I}^P\}$ 
4   do  $T_{max} \leftarrow 0$ 
5      $\mathbf{w}_{max} \leftarrow \text{none}$ 
6     for each  $\mathbf{w}_j^a$  in  $templates$ 
7       do if  $[\rho(\mathbf{I}^r, \mathbf{w}_j^a) \geq \rho_a] \ \& \ [T(\mathbf{I}^r, \mathbf{w}_j^a, \beta_a) > T_{max}]$ 
8         then
9            $T_{max} \leftarrow T(\mathbf{I}^r, \mathbf{w}_j^a, \beta_a)$ 
10           $\mathbf{w}_{max}^a \leftarrow \mathbf{w}_j^a$ 
11          if  $\mathbf{w}_{max}^a \neq \mathbf{w}_0 \ \& \ class(\mathbf{I}^r) = class(\mathbf{w}_{max}^a)$ 
12            then  $\mathbf{w}_{j_{max}} \leftarrow \mathbf{w}_{j_{max}} \wedge \mathbf{I}^r$ 
13            else  $templates \leftarrow templates \cup \{\mathbf{I}^r\}$ 
14 return  $templates$ 

```

Fig. 4. Anagnostopoulos' no-matchtracking FAM.

When dividing the network we must guarantee that certain consistency properties hold. These properties have been formally proven and are simply stated below as theorems.

Non-Duplication. *A template \mathbf{w} will either be owned by a single processor, or it will be in transit on a single processor (i.e. only one copy of the template exists in the system).*

Bundle size. *The excess of templates for a process $k \neq 0$, at any given time, always fits in the packet size to be sent back to the pipeline.*

No overflow. *The first processor in the pipeline can always absorb all templates that are sent back to it.*

Bounded delay. *The templates that a process k has, on the current iteration, were created at least $p - k - 1$ iterations ago where p is the total number of processors in the pipeline.*

Pipeline depth invariance. *The difference in the number of templates that 2 arbitrary processes in the pipeline have cannot exceed $p_a + 1$, where p_a is the packet size.*

4.3. Network partitioning versus data partitioning

One of the advantages of data partitioning for a quadratic complexity algorithm is that by splitting the problem into independent subsets the convergence time is reduced proportionally to the square of the number of partitions. Unfortunately, this approach alters the nature of the algorithm. Therefore, when data partitioning is applied to FAM it is important to demonstrate that the classification performance of the partitioned FAM is comparable to the classification of the original FAM. On the other hand, network partitioning using a pipeline is identical to the original no-matchtracking FAM. Nevertheless, pipelined FAM can only speed up FAM's complexity linearly with respect to the number of processors in the pipeline.

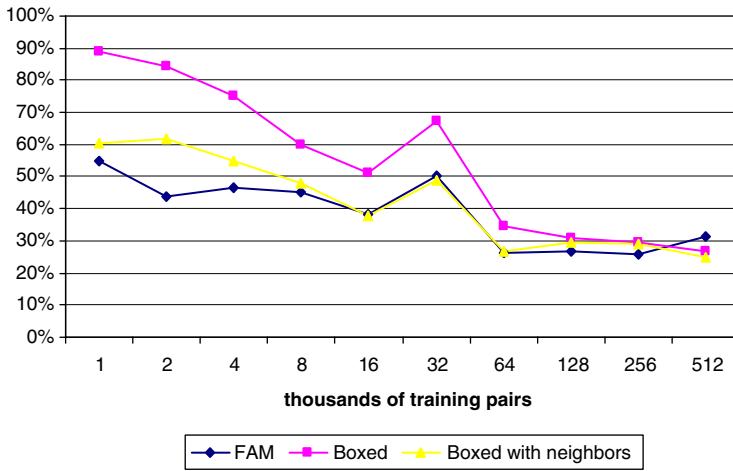


Fig. 5. Forest Covertypes database error rate versus training set size (in thousands of patterns).

5. Experimental results

To prove the feasibility of our approaches we conducted a series of tests on the Covertypes database by Blackard [9]. The size of the data used for training started from 1000 data-points and increased incrementally, by a factor of 2, until it reached a training set size of 512,000 data-points. Classification performance was evaluated using a fixed set of 20,000 test patterns (see Fig. 5).

5.1. Partitioned FAM

Comparisons of the training performance of the partitioned FAM approach and the non-partitioned FAM approach on the aforementioned database were conducted. The performance comparison was based on two measures: the time that it took for the networks to undergo one epoch of training (see results in Fig. 6) and the generalization performance of the trained networks on the chosen test sets (see Fig. 5). The dimensions to project the data in the partitioned-FAM approach were chosen manually by simply observing the range and variation of the values of the datasets across the chosen dimensions. A vigilance value of 0.96 was used in the training of FAM. This gives a partitioning scheme with boxes of side size of 0.04, resulting in $25^3 = 15,625$ boxes in the three dimensions that we used to project the data.

Some of the observations from the results (error rate and training time) of the FAM and partitioned FAM with these databases are: (a) the partitioned FAM reduces the training time compared to FAM, and at times significantly (especially when the size of the training set is large) (b) the generalization performance of the partitioned FAM is slightly inferior to FAM, especially for training sets of smaller size.

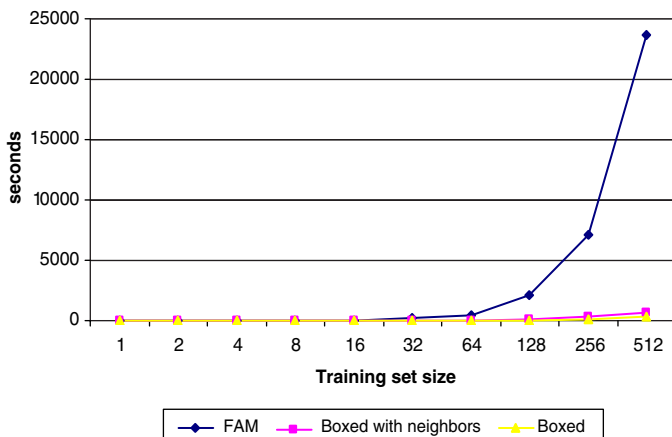


Fig. 6. Forest Coverttype database epoch training time (in seconds) against training set size (in thousands of patterns).

Table 1
Training Time of the Pipelined no-matchtracking, FAM implemented on a Beowulf cluster of workstations

Proc	DBs:32	DBs:64	DBs:128	DBs:256	DBs:512
1	9.3	28.8	96.0	357.0	1367.9
2	4.8	14.4	49.9	155.4	552.7
4	2.5	7.9	22.1	81.6	273.2
8	1.4	4.2	11.0	51.5	140.2
16	0.9	2.0	7.6	23.4	76.5
32	0.9	1.6	4.5	14.6	38.3

The processor 1 case corresponds to implementation on a sequential machine, DBs stands for database size in thousands of patterns.

It is worth mentioning that the name “boxed” in Figs. 5 and 6 refers to the partitioned-FAM approach where the trained FAMs take into consideration only the data-points residing in their corresponding box. The term “boxed with neighbors” is the partitioned FAM approach where the trained FAMs take into consideration not only the data-points in their corresponding box but also the data-points in the boxes neighboring their corresponding box. Obviously, the term “FAM” corresponds to a single FAM, trained with all the data-points in the training set.

5.2. Pipelined FAM

Classification results were comparable to those of the original FAM algorithm when using the pipelined no-matchtracking FAM on the Coverttype database. On the other hand, speedup was substantial, and as expected linear with respect to the number of processors used in the pipeline (see Table 1).

6. Conclusions

Two partitioning schemes were proposed for the FAM algorithm to work on a parallel setting. Network-partitioned FAM exhibited linear speedup with respect to the number of processors, while data-partitioned FAM provided quadratic speedup with respect to the number of processors used. These results show promise and encourage us to continue with further research. Among our current interests we mention: (a) the development of a pipelined matchtracking FAM parallel algorithm (b) the addition of a pruning scheme for the Boxing approach to improve its compression ratio and (c) the use of other parallel designs like master–slave arrangements that require broadcast send and broadcast receive operations.

Acknowledgements

The authors would like to thank the Institute of Simulation and Training and the Link Foundation Fellowship program for partially funding this project. This work was also supported in part by the National Science Foundation under Grant no. 0203446.

References

- [1] G.C. Anagnostopoulos, Putting the utility of match tracking in Fuzzy ARTMAP to the test, in: *Proceedings of the Seventh International Conference on Knowledge-Based Intelligent Information Engineering*, vol. 2, University of Oxford, UK, 2003, pp. 1–6, KES'03.
- [2] L. Breiman, J.H. Friedman, R.A. Olshen, J. Stone, *Classification and Regression Trees*, CRC Press, Boca Raton, FL, 1998.
- [3] G.A. Carpenter, S. Grossberg, N. Markuzon, J.H. Reynolds, D.B. Rosen, Fuzzy ARTMAP: A neural network architecture for incremental learning of analog multidimensional maps, *IEEE Trans. Neural Networks* 3 (5) (1992) 698–713.
- [4] G.A. Carpenter, S. Grossberg, J.H. Reynolds, Fuzzy ART: An adaptive resonance algorithm for rapid, stable classification of analog patterns, in: *International Joint Conference on Neural Networks, IJCNN'91*, vol. II, IEEE/INNS/ENNS, Seattle, Washington, 1991, pp. 411–416.
- [5] M. Mehta, R. Agrawal, J. Riassnen, SLIQ: A fast scalable classifier for data mining, in: P.M.G. Apers, M. Bouzeghoub, G. Gardarin, *Extending Database Technology*, Avignon, France, Springer, Berlin, 1996, pp. 18–32. [Online]. Available: <http://citeseer.mj.nec.com/mehta96sliq.html>.
- [6] C.M. Procopiuc, M. Jones, R. Agrawal, T. Murali, A Monte Carlo algorithm for fast projective clustering, ser. ACM SIGMOD, ACM, Ed., Madison, Wisconsin, USA, June 2002, pp. 418–427.
- [7] J.R. Quinlan, *C4.5: Programs for Machine Learning*, Morgan Kaufmann, San Mateo, California, 1993.
- [8] J.C. Shafer, R. Agrawal, M. Mehta, SPRINT: A scalable parallel classifier for data mining, in: *Proceedings of the 22nd International Conference on Very Large Databases, VLDB*, T.M. Vijayaraman, A.P. Buchmann, C. Mohan, N.L. Sarda (Eds.), Bombay, India, Morgan Kaufmann, September 1996, pp. 544–555.
- [9] University of California, Irvine, UCI machine learning repository, <http://www.ics.ucf.edu/mllearn/MLRepository.html>, 2003.