

**NEW CLASSES OF LEARNING AUTOMATA BASED SCHEMES
FOR ADAPTATION OF BACKPROPAGATION
ALGORITHM PARAMETERS ***

M. R. MEYBODI ** AND H. BEIGY

Dept. of Computer Engr., Amirkabir University of Technology, Tehran, I. R. of Iran

Abstract – One popular learning algorithm for feedforward neural networks is the back-propagation (BP) algorithm which includes parameters: learning rate (η), momentum factor (α) and steepness parameter (λ). The appropriate selections of these parameters have a large effect on the convergence of the algorithm. Many techniques that adaptively adjust these parameters have been developed to increase speed of convergence. In this paper, we shall present several classes of learning automata based solutions to the problem of adaptation of BP algorithm parameters. By interconnection of learning automata to the feedforward neural networks, we use learning automata schemes for adjusting the parameters η , α , and λ based on the observation of random response of the neural networks. One of the important aspects of proposed scheme is its ability to escape from local minima with high possibility during the training period. The feasibility of the proposed methods are shown through the simulations on several problems.

Keywords – Neural network, back-propagation, fixed structure learning automata, steepness parameter

1. INTRODUCTION

Back-propagation algorithm is a systematic method for training multilayer neural networks. Despite the many successful applications of back-propagation, it has many drawbacks. For complex problems it may require a long time to train the networks, and it may not train at all. Long training time can be the result of the non-optimum values for the parameters of the training algorithm. It is not easy to choose appropriate values for these parameters for a particular problem. The parameters are usually determined by trial and error and using past experiences. For example, if the learning rate is too small, convergence can be very slow, if too large, paralysis and continuous instability may result. Moreover, the best value at the beginning of training may not be so good later. Thus several researches have suggested algorithms for automatically adjusting the parameters of training algorithm as training proceeds.

Arabshahi *et al.* [1] proposed an error back-propagation algorithm in which the learning-rate is adapted. In this algorithm, learning-rate is a function of error and changes in the error. They proposed that the learning-rate be adjusted using a fuzzy logic control system, in which the error and changes in error are the inputs and changes in learning-rate in the output of fuzzy logic controller. Kandil *et al.* [2] used optimum, time-varying learning-rate for multi-layer neural network by linearizing the neural

*Received by the editors June 22, 1998 and in final revised form December 3, 2000

**Corresponding author

network around weight vector at each iteration. Parlos *et al.* [3] proposed an accelerated learning algorithm for supervised training of multi-layer neural networks named Adaptive Error Back-Propagation Algorithm. In their proposed algorithm the learning-rate is a function of the error and the error gradient. Ref. [4-13], have proposed other schemes for adaptation of learning rate. Sperduti and Starita [14] proposed an error back-propagation algorithm in which the steepness parameter is adapted using gradient descent algorithm. Several learning automata (LA) based procedures have also been developed [15-20]. In these methods variable structure learning automata (VSLA) or fixed structure learning automata (FSLA) have been used to find the appropriate values of parameters for the BP algorithm. In these schemes either a separate learning automata is associated to each layer of the network or a single automata is associated to the whole network to adapt the appropriate parameters. It is shown that the learning rate adapted in such a way, not only increases the rate of convergence of the network, but it by-passes the local minimum in most cases.

In this paper, we propose two new classes of LA based schemes for adaptation of appropriate learning rate or steepness parameters for BP algorithm. Unlike the existing LA based schemes, in these schemes one learning automata is assigned to every link or every neuron in the network for determining the parameters for that link or neuron. The simulation results show the feasibility of the proposed method and its superiority to the existing LA based schemes. The proposed schemes have two important aspects: higher speed of convergence and a higher probability of escaping from the local minima. In order to evaluate the performance of proposed schemes, simulations are carried out on four learning problems: digit recognition, encoding, odd parity, and symmetry problems and the results are compared with results obtained from standard BP.

The rest of the paper is organized as follows: Section 2 briefly presents the basic back-propagation algorithm and learning automata. Application of learning automata for adaptation of learning rate, momentum factor, and steepness parameter is given in Section 3. Section 4 presents the proposed learning automata based schemes. The simulation results are given in Section 5. Section 6 discusses the time and space complexity of the proposed method. Section 7 concludes the paper.

2. BACKPROPAGATION ALGORITHM AND LEARNING AUTOMATA

In this section, in all brevity, we discuss the fundamentals of back-propagation algorithm and learning automata.

a) Back-propagation algorithm

Error backpropagation algorithm, which is an iterative gradient descent algorithm, is a simple way to train multilayer feedforward neural networks [21]. The BP algorithm is based on the gradient descent rule:

$$W(n+1) = W(n) + \eta G(n) + \alpha [W(n) - W(n-1)] \quad (1)$$

where, W is the weight vector, n is the iteration number, η is learning rate, α is momentum factor, and G is gradient of error function that is given by:

$$G(n) = -\nabla E_p(n) \quad (2)$$

where, E_p is the sum of squared error given by:

$$E_p(n) = \frac{1}{2} \sum_{j=1}^{\#outputs} [T_{p,j} - O_{p,j}]^2 \quad (3)$$

for $p=1,2,\dots, \#$ patterns. Where, $T_{p,j}$ and $O_{p,j}$ are desired and actual outputs for pattern p at output node j . One of the major problems encountered during implementation of the BP algorithm is proper choice and update of the learning rate η to allow convergence, while keeping the number required iterations at a reasonable number. One of the main reasons for investigating the possibility of the adaptive learning rate rule, is the desire to reduce the sensitivity of the learning on the learning rate, without adding more tuning parameters.

In the BP algorithm framework, each computational unit computes the same activation function. The computation of the sensitivity for each neuron requires the derivative of activation function, therefore, this function must be continuous and differentiable. The activation function is normally a sigmoid function chosen between $1/(1+\exp(-\lambda net))$ and $\tanh(\lambda net)$. The coefficient of the exponent of the exponential term determines the steepness of linearity of that function. The steepness parameter λ is often set to a constant value and not changed by the learning algorithm. We gain much flexibility if we move the net inputs of the sigmoidal functions near to their active regions, where the associated gradient are not very close to zero. This makes the BP algorithm not trapped at some points in the network parameters space where the BP algorithm would effectively stop, even though it is not close to a local minima point. This will cause the gradient of the error function to be small if the sigmoidal is shifted far outside the active region of the input to the function. Therefore, it is better to center each sigmoid to be inside the active region of the sigmoidal function.

The momentum term in weight adaptation Eq. (1) causes a large change in the weight if the changes are currently large, and will decrease as the changes become less. This means that the network is less likely to get stuck in local minima early on, since the momentum term will push the changes over local downward trend. Momentum is of great assistance in speeding up convergence along shallow gradients, allowing the path the network takes toward the solution to pickup speed in the downhill direction. The error surface may consist of long gradually sloping ravines which finish at minima. Convergence along these ravines is slow, and usually the algorithm oscillates across the ravine valley as it moves towards a solution. This is difficult to speed up without increasing the chance of overshooting the minima, but the addition of the momentum term is fairly successful. This difficulty could be removed if we select the momentum factor to be small at the near of minima and large far from minima.

b) Learning automata

Learning automata (LA) can be classified into two main families, fixed and variable structure learning automata [22, 25]. Examples of the FSLA type that we use in this paper are Tsetline, Krinsky, Tsetline G, and Krylov automata. A fixed structure learning automaton is a quintuple $\langle \alpha, \Phi, \beta, F, G \rangle$ where:

- 1) $\alpha = (\alpha_1, \dots, \alpha_R)$ is the set of actions that it must choose from.
- 2) $\Phi = (\Phi_1, \dots, \Phi_S)$ is the set of states.
- 3) $\beta = \{0,1\}$ is the set of inputs where 1 represents a penalty and 0 represents a reward.
- 4) $F : \Phi \times \beta \rightarrow \Phi$ is a map called the transition map. It defines the transition of the state of the automation on receiving an input F may be stochastic.
- 5) $G : \Phi \rightarrow \alpha$ is the output map and determines the action taken by the automaton if it is in state Φ_j .

The selected action serves as the input to the environment which in turn emits a stochastic response $\beta(n)$ at the time n . $\beta(n)$ is an element of $\beta = \{0, 1\}$ and is the feedback response of the environment to the automaton. The environment penalizes (i.e., $\beta(n) = 1$) the automaton with the penalty c_j , which is the action dependent. On the basis of the response $\beta(n)$, the state of the automaton $\Phi(n)$ is updated and a new action chosen at the time $(n+1)$. Note that the $\{c_i\}$ are unknown initially and it is desired that as a result of interaction with the environment, the automaton arrive at the action, which presents it with the minimum penalty response, in an expected sense. If the

probability of the transition from one state to another state and probabilities of correspondence of action and state are fixed, the automaton is said fixed-structure automata and otherwise the automaton is said variable-structure automata. We summarize some fixed-structure learning automaton and variable structure automaton in the following paragraphs.

c) *The two-state automata ($L_{2,2}$)*

This automata has two states, ϕ_1 and ϕ_2 and two actions α_1 and α_2 . The automata accepts input from a set of $\{0, 1\}$ and switches its states upon encountering an input 1 (unfavorable response) and remains in the same state on receiving an input 0 (favorable response). An automata that uses this strategy is referred to as $L_{2,2}$ where the first subscript refers to the number of states and second subscript to the number of actions.

d) *The two-action automata with memory ($L_{2N,2}$)*

This automata has $2N$ states and two actions and attempts to incorporate the past behavior of the system in its decision rule for choosing the sequence of actions. While the automata $L_{2,2}$ switches from one action to another on receiving a failure response from environment, $L_{2N,2}$ keeps an account of the number of success and failures received for each action. It is only when the number of failures exceeds the number of successes, or some maximum value N , the automata switches from one action to another. The procedure described above is one convenient method of keeping track of performance of the actions α_1 and α_2 . As such, N is called memory depth associated with each action, and automata is said to have a total memory of $2N$. For every favorable response, the state of automata moves deeper into the memory of corresponding action, and for an unfavorable response, moves out of it. This automaton can be extended to multiple action automata and this automaton is named $L_{2N,2}$ automata. The state transition graph of $L_{2N,2}$ automata is shown in Fig. 1

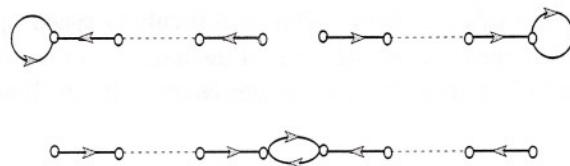


Fig. 1. The state transition graph for $L_{2N,2}$

e) *The Krinsky automata*

This automata behaves exactly like $L_{2N,2}$ automata when the response of the environment is unfavorable, but for favorable response, any state ϕ_i (for $i = 1, \dots, N$) passes to the state ϕ_1 and any state ϕ_i (for $i = N+1, \dots, 2N$) passes to the state ϕ_{N+1} . This implies that a string of N consecutive unfavorable responses are needed to change from one action to another. The state transition graph of Krinsky automata is shown in Fig. 2

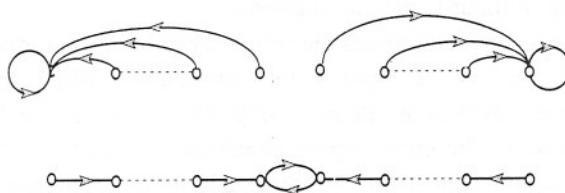


Fig. 2. The state transition graph for Krinsky Automata

f) The Krylov automata

This automata has state transition that are identical to the $L_{2N,2}$ automata when the output of the environment is favorable. However, when the response of the environment is unfavorable, a state ϕ_i ($i \neq 1, N, N+1, 2N$) pass to a state ϕ_{i+1} with probability 0.5 and to a state ϕ_{i-1} with probability 0.5. When $i = 1$ or $i = N+1$, ϕ_i stays in the same state with probability 0.5 and moves to ϕ_{i+1} with the same probability. When $i = N$, automata state moves to state ϕ_{N-1} and ϕ_{2N} with the same probability 0.5. When $i=2N$, automata state moves to state ϕ_{2N-1} and ϕ_N with the same probability 0.5. The state transition graph of Krylov automata is shown in Fig. 3.

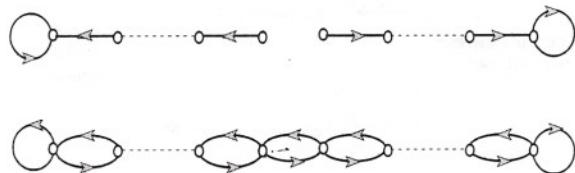


Fig. 3. The state transition graph for Krylov Automata

In this paper we refer to an automata by the name of automata followed by the list of parameters for that automata, the first parameter refers to the number of actions and the second parameter refers to the depth for each action.

g) Variable structure automata

Variable-structure automata is represented by sextuple $\langle \underline{\beta}, \underline{\phi}, \underline{\alpha}, \underline{P}, G, T \rangle$, where, $\underline{\beta}$ is a set of inputs actions, $\underline{\phi}$ is a set of internal states, $\underline{\alpha}$ a set of outputs, \underline{P} denotes the state probability vector governing the choice of the state at each stage k , G is the output mapping, and T is learning algorithm. The learning algorithm is a recurrence relation and is used to modify the state probability vector [22].

It is evident that the crucial factor affecting the performance of the variable structure learning automata, is learning algorithm for updating the action probabilities. Various learning algorithms have been reported in the literature [22]. Let α_i be the action chosen at time k as a sample realization from distribution $p(k)$. The linear reward-penalty algorithm (L_{R-P}) is one of the earliest schemes. In an L_{R-P} scheme the recurrence equation for updating p is defined as

$$\beta(k) = 0 \Rightarrow p_j(k) = \begin{cases} p_j(k) + a(l - p_j(k)) & \text{if } i = j \\ p_j(k)(1 - a) & \text{if } i \neq j \end{cases}$$

$$\beta(k) = 1 \Rightarrow p_j(k) = \begin{cases} p_j(k)(l - b) & \text{if } i = j \\ \frac{b}{r - l} + (l - b)p_j(k) & \text{if } i \neq j \end{cases}$$

The parameters a and b represent reward and penalty parameters, respectively. The parameter $a(b)$ determines the amount of increase (decreases) of the action probabilities.

3. LA BASED SCHEMES FOR ADAPTATION OF BP PARAMETERS

In this section, we first, briefly describe previous LA based schemes [15-20] for adaptation of BP parameters and then introduce two new classes of LA based schemes. In all of the existing schemes,

one or more automata have been associated to the network. The learning automata based on the observation of the random response of the neural network, adapt one or more of BP parameters. The interconnection of learning automata and neural network is shown in Fig. 4. Note that the neural network is the environment for the learning automata. The learning automata according to the amount of the error received from neural network adjust the parameters of the BP algorithm. The actions of the automata correspond to the values of the parameters being calculated and input to the automata is some function of the error in the output of neural network.

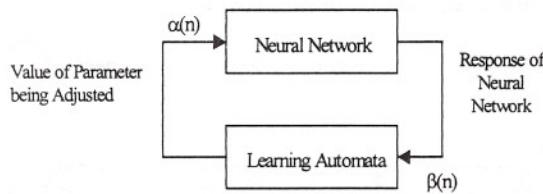


Fig. 4. The connection of learning automata and neural network

A function of error between the desired and actual outputs of network is considered as the response of environment. Windows on the past values of the errors are swiped and the average value of the error in this window computed. If the difference of the average value in the two last steps is less than the predefined threshold value, the response of the environment is favorable and if this difference of average value is greater than the threshold value, the response of the environment is unfavorable.

Existing LA based procedures for adaptation of BP parameters can be classified into two groups, which we refer to as, group **A** and group **B**. In group **A** schemes, an automaton is used for the whole network [15,17,20] whereas in group **B** schemes, separate automata one for each layer (hidden and output) are used [16, 18-20]. Each group **A** and **B** depending on the type of automata used (fixed or variable structure) can be classified into two subgroups. The parameter adapted by group **A** schemes will be used by all the links or neurons of the networks and, therefore, these schemes fall into the category of global parameter adaptation method Group **B** schemes by adapting the parameter for each layer independently, may be referred to as quasi-global parameter adaptation methods. For the sake of convenience in presentation, we use the following naming conventions to refer to different LA based schemes in class **A** and class **B**.

Automata-AV(γ): A scheme in class **A** for adjusting parameter γ , which uses variable structure learning automata **Automata**.

Automata-AF(γ): A scheme in class **A** for adjusting parameter γ , which uses fixed structure learning automata **Automata**.

Automata₁-Automata₂-BV(γ): A scheme in class **B** for adjusting parameter γ , which uses variable structure learning automata **Automata₁** for hidden layer and variable structure learning automata **Automata₂** for output layer in a three layers network.

Automata₁-Automata₂-BF(γ): A scheme in class **B** for adjusting parameter γ which uses fixed structure learning automata **Automata₁** for hidden layer and fixed structure learning automata **Automata₂** for output layer in a three layers network.

The letters F and V in the above names denotes FSLA and VSLA, respectively. For all the LA based schemes reported, it is shown through simulation that the use of LA for adaptation of BP learning algorithm parameters increases the rate of convergence by a large amount. Figure 5 borrowed from [20], compares the effectiveness of different LA based schemes in class **A** for adaptation of learning rate for the 8×8 dots numeric font recognition problem. In this simulation, the threshold of 0.01 and window size of 1 is chosen. For linear reward-penalty automata, the reward and penalty coefficient 0.001 and 0.0001 are chosen. It is reported that FSLA based schemes perform

better than the VSLA based schemes [17]. Also, simulation studies have shown that by using LA based schemes for adaptation of learning rate or momentum factor, we can compute a new point that is closer to the optimum than the point computed by BP algorithm which, uses a fixed pre-determined learning rate or momentum factor [18-19].

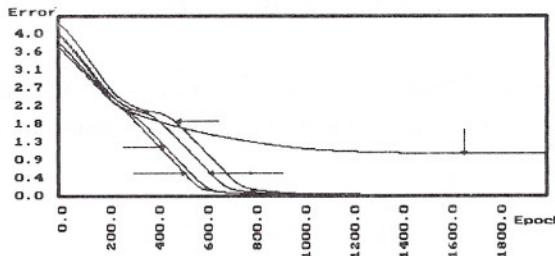


Fig. 5. Performance of different class A based schemes a: Standard BP b: Tsetline(4,4)-AF (η)
c: Krinsky(2,4)-AF (η) d: Krylov(2,4)-AF (η) e: L_{RP}-AV (η)

4. NEW CLASSES OF SCHEMES FOR ADAPTATION OF BP PARAMETERS

In this section, we propose two new classes of LA based schemes called classes **C** and **D**. In a class **C** scheme one automata is associated to each link of the network to adjust the parameter for that link and in a class **D** scheme, one automata is associated to each neuron of the network to adjust the parameter for that neuron. Group **C** and **D** schemes may be referred to as the local parameter adaptation methods. We use **Automata-CV(γ)** and **Automata-CF(γ)** to refer to the schemes in class **C** and **Automata-DV(γ)** and **Automata-DF(γ)** to refer to the schemes in class **D**. The letters F and V in above names denote FSLA and VSLA, respectively. We use class **C** schemes for adaptation of learning rate and/or momentum factor and class **D** schemes for adaptation of steepness parameter. In class **C** and **D** schemes, the automata receives favorable response from the environment if the algebraic sign of derivative in two consecutive iterations is the same and receives unfavorable response if the algebraic sign of the derivative in two consecutive iterations alternates. The following algorithms describe class **C** and class **D** schemes.

```

procedure C_Scheme_BP (Automata)
    Initialize the weights to small random values.
    initialize the parameters for automaton Automata.
repeat
    for all training patterns (X, T) in the training set do
        Call FeedForward
        Call ComputeGradient
        for all layers in the network do
            for all nodes in lth layer do
                for all weights w for nth node in lth layer do
                    if Sign ( $\frac{\partial E_p}{\partial w}$ , k) = Sign ( $\frac{\partial E_p}{\partial w}$ , k-1) Then
                        //The sign at iteration k and k-1 is the same
                         $\eta$  = Call Automata (0) //  $\beta$  is 0
                    else
                         $\eta$  = Call Automata (1) //  $\beta$  is 1
                    End if
                End for
            End for
        End for
    End for
    Call UpdateWeights // Batch weights updating is used
    until k > N. // N is maximum training epoch
End procedure

```

Fig. 6. BP with a class **C** scheme

```

Initialize the weights to small random values.
initialize the parameters for the automaton Automata.
repeat
    for all training patterns (X, T) in the training set do
        Call FeedForward
        Call ComputeGradient
        for all layers in the network do
            for all steepness parameters  $\lambda$  in lth layer do
                if Sign ( $\frac{\partial E_p(k)}{\partial \lambda}$ ) = Sign ( $\frac{\partial E_p(k-1)}{\partial \lambda}$ ) Then
                    //The sign at iteration k and k-1 is the same
                     $\lambda$  = Call Automata (0) //  $\beta$  is 0
                else
                     $\lambda$  = Call Automata (1) //  $\beta$  is 1
                End if
            End for
        End for
        Call UpdateWeights
        // Batch weights updating is used
    until k > N // N is maximum training epoch
End procedure

```

Fig. 7. BP with a class **D** scheme

In classes **D** schemes, to compute the partial derivative of error with respect to steepness parameter ($\partial E_p / \partial \lambda$), we minimize the error given by Eq. (3). Assume that net_I^L represents the net output of I^{th} neuron in L^{th} layer, which is given by

$$net_I^L = \sum_{k=0}^{N_{L-1}} W_{K,I}^L \times O_K^{L-1}$$

O_K^L is output of k^{th} neuron in L^{th} layer given by $O_K^L = f_K^L(net_K^L)$, where f_K^L shows the activation function of k^{th} neuron in L^{th} layer which is one of sigmoidal or tanh functions. Differentiation of E with respect to λ_K^L yields to:

$$\frac{\partial E_p}{\partial \lambda_K^L} = -\delta_K^L \times \frac{\partial f_K^L}{\partial \lambda_K^L}$$

where

$$\delta_K^L = \begin{cases} (T_K - O_K) & \text{if } L \text{ is output layer} \\ \sum_{J=1}^{N_{L+1}} \delta_J^{L+1} \times W_{K,J}^{L+1} \times \frac{\partial f_K^L}{\partial net_J^{L+1}} & \text{if } L \text{ is hidden layer} \end{cases}$$

a) A deterministic fixed structure learning automaton

In what follows, we introduce a new fixed structure learning automata called J-automata. This automaton can be used by class **C** and class **D** schemes to achieve higher degree of performance than the other known schemes in class **C** or class **D**. The proposed automata which we denote by J_{KNK} has KN states and K actions and attempts to incorporate the past behavior of the system in its decision rule for choosing the sequence of actions. States with numbers $(k-1)N+1$ through kN correspond to action k . The state transition graph of this automaton for favorable response and unfavorable response is shown in Fig. 8.

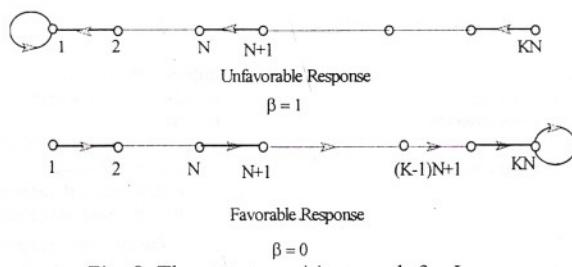


Fig. 8. The state transition graph for $J_{KN,K}$

The environment produces the favorable response if the algebraic sign of the derivative in two consecutive iterations is the same, and unfavorable response produced by the environment if the algebraic sign of the derivative in two consecutive iterations alternates. When this automaton is used to adjust the learning rate (steepness parameter), each action corresponds to one of the values of the learning rate (steepness parameter). The automata, on reward, move to the higher number state (hoping eventually it increases the value of the parameter) and on penalty move to the lower number state (hoping eventually it decreases the value of the parameter). It is only when the sign of derivative remains the same for N consecutive iterations or alternates for N consecutive iteration the automata changes its action (switches from one value for the related parameter to another value). Therefore, N is memory depth associated with each action and automata is said to have a total memory of KN . If

parameter μ can assume K values, $\mu_1 \leq \mu_2 \leq \dots \leq \mu_K$ then states $(k-1)N + 1$ through kN correspond to the value of μ_k . Clearly such an assignment of values of parameter μ to the states of the automata causes the value used by BP to increase, if in N consecutive iterations, the sign of the derivative do not change and decreases, if in N consecutive iterations the sign of derivative alternates.

Remark 1: In this remark we first explain the variable learning rate scheme (VLR) and then explain why $J_CF(\eta)$ performs better than the VLR scheme. Variable learning rate is a scheme in which the learning rate is varied according to the performance of the algorithm [13]. If the error decreases after a weight update, then the learning rate is increased by some factor (e.g., 1.05). If the error increases more than some set of percentage (typically one to five percent), then the weight update is discarded and the learning rate is decreased by some factor (e.g., 0.7) and the momentum term (if it is used) is set to zero. When a successful step is taken, the momentum term is reset to its original value. If the algorithm is working well, and the error continues to go down, then learning rate will increase and convergence will speedup.

Assume that the network has n adjustable parameters π_1, \dots, π_n . Gradient of error function with respect to π_1, \dots, π_n is defined as:

$$\nabla E = \left[\frac{\partial E}{\partial \pi_1}, \frac{\partial E}{\partial \pi_2}, \dots, \frac{\partial E}{\partial \pi_n} \right]^T$$

Gradient vector ∇E extends in the direction of the greatest rate of change of E , but it doesn't mean that all $\nabla_k E$ (the k^{th} element of ∇E) have the same algebraic sign. In other words, the greatest rate of increase (decrease) of E doesn't imply the same sign for projections of $\nabla_k E$ along all axis.

VLR scheme decreases (increases) learning rate as the result of increase (decrease) of E based on gradient. In the VLR scheme a single learning rate will be adapted and be used by BP to walk along all directions, which may create oscillation along some axis and leads to lower rate of convergence. On the contrary, in **C** (**D**) schemes, since each link (neuron) has its own learning rate (steepness) (that is, the learning rate (steepness parameter) for every axis is adapted independently), oscillation may be decreased, and as a result, a higher rate of convergence can be obtained.

b) Simultaneous adaptation of learning rate and steepness parameters

The rate of convergence can be improved if, both learning rate and steepness parameter are adapted simultaneously. For simultaneous adaptation of learning rate and steepness parameters, a class **C** scheme is used for adaptation of learning rate and a class **D** scheme for adaptation of steepness parameter. A scheme that simultaneously adapts learning rate and steepness parameter is denoted by **Automata1-Automata2-CDF**(μ, λ), if FSLA is used and **Automata1-Automata2-CDV**(μ, λ), if VSLA is used. The **Automata1** is used for adaptation of learning rate and **Automata2** for adaptation of steepness parameter.

c) Simultaneous adaptation of learning rate and momentum factor

A simple method of increasing the speed of learning and stability of training algorithm is to modify the standard BP by including the momentum factor [21] which is the simplest form of BP with momentum term given in Eq. (1). The changes in the weight at iteration of n is given by:

$$\Delta W(n) = \eta G(n) + \alpha \Delta W(n-1) \quad (4)$$

Solving the difference equation (4) gives to the following time series equation.

$$\Delta W(n) = \eta \sum_{t=0}^n \alpha^{n-t} G(t) \quad (5)$$

- By inspection of Eq. (5), we may make the following useful observations:
The current adjustment of $\Delta W(n)$ represents the sum of an exponentially weighted time series. This equation converged if, and only if, $0 \leq |\alpha| < 1$.
- When G has the same algebraic sign on consecutive iterations, $|\Delta W|$ grows, and W is adjusted by a large amount. Hence, the inclusion of momentum term accelerates the convergence of algorithm. To accelerate more, the momentum factor must have a large value as much as possible.
- When the algebraic sign of G alternates in consecutive iterations, the $|\Delta W|$ shrinks and W is adjusted by a small amount. Hence, the inclusion of momentum term stabilizes the convergence of algorithm. To accelerate more, the momentum factor must have a small value as much as possible.

From the above observation, we may conclude that the momentum factor could be adjusted by $J_{KN,K}$ automata in the same manner as the learning rate. We denote a scheme that simultaneously adapted both learning rate and momentum factor by **Automata1-Automata2-CF** (η, α). The **Automata1** is used for adaptation of learning rate and **Automata2** for adaptation of steepness parameter.

Learning rate can be adjusted directly or indirectly. Momentum, conjugate gradient, and class **D** are examples of indirect methods. Since the indirect methods have not been satisfactory enough in many cases, the direct adjustment of learning rate have been proposed. Direct methods can be classified in three groups: global methods, such as Bold Driver and Class **A** schemes, in which one learning rate is used for all weights, quasi-global methods, such as class **B** schemes, in which one learning rate is used for a set of weights, and local methods in which one learning rate is used for each weight in the network. Examples of local method are SAB, SuperSAB, and class **C** schemes.

With the introduction of these new classes of schemes we propose a new classification tree for learning rate adjustment methods. Figure 9 shows this new classification.

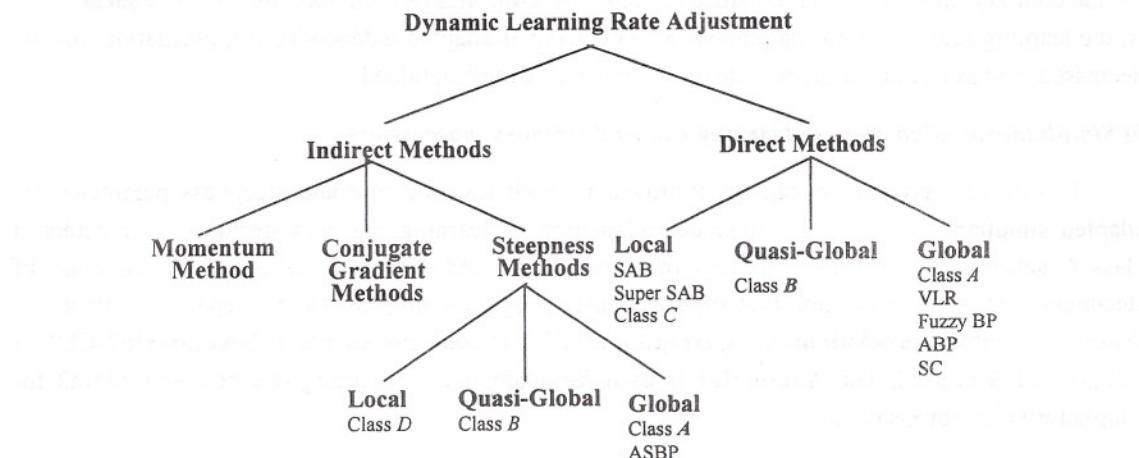


Fig. 9. Classification tree for dynamic learning rate adjustment schemes

5. SIMULATIONS

In order to evaluate the performance of proposed schemes, simulations are carried out on four learning problems: digit recognition, encoding, odd parity, and symmetry problems and the results are compared with results obtained from standard BP. These problems are chosen because they possess different error surfaces and collectively present an environment that is suitable to determine the effect

of the proposed method. In this section, we first describe the test problems and then present the results of the simulation.

a) 8×8 Dot numeric font learning

There are numbers 0, ..., 9, and each represented by a 8×8 grid of black and white dots as shown in Fig. 10 [14]. The network must learn to distinguish these numbers. The network architecture used for this problem, consists of 64 input units, which are connected to 6 hidden units, which are connected to 10 output units.

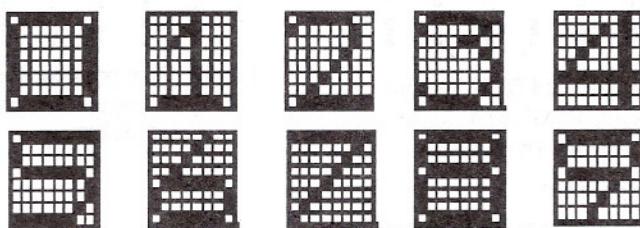


Fig. 10. Grid of black and white dots

b) Three-bit odd-parity problem

In this problem, a string of three inputs is applied to the network, the output of network is zero (one) if number of ones in the input is odd (even) [21]. The training set consists of 8 patterns. This network has three input units, three hidden units, and one output unit.

c) Encoding problem

In this problem, a set of orthogonal input patterns are mapped to a set of orthogonal output patterns through a small set of hidden units [21]. The training set consists of 8 patterns. The network architecture used for solving the problem consists of 8 input units, 3 hidden units, and 8 output units.

d) Symmetry problem

This problem classifies input strings, as to whether or not they are symmetric about center [21]. The training set consists of 64 patterns. The network architecture used for solving this problem consists of 8 input units, 2 hidden units, and 1 output unit.

Typical simulations for these four problems for different parameter adaptation schemes are shown in Figs. 11 through 25. Figure 11 compares the performance of different class **A** schemes with J_CF scheme. Figure 12 indicates that J_CF scheme has a higher speed of convergence than any known scheme in class **C**. Figure 13 compares the J_CF scheme with the best scheme in class **B**. Figures 14 through 16 indicate that, if η and λ adapted simultaneously, the performance of the BP algorithm increases by a large amount. Figures 17 through 19 compare the performance of J_J_CDF (η, λ) with VLR scheme and schemes proposed by Arabshahi (Fuzzy BP) and Darken and Moody (SC). The SAB and SuperSAB schemes are tested on parity, encoding, and numeric font recognition problems and compared with standard BP, VLR, J_CF, and J_CDF schemes. The simulation results, which are given in Figs. 20 through 22, show the effectiveness of J_CF and J_CDF schemes. The Qprop and RPROP methods [9] are simulated on parity and encoding problems and compared with standard BP, J_CF, J_DF, and J_CDF schemes. The simulation results show that in some cases, the RPROP exhibits the same performance as the proposed schemes, show the superiority of J_CF and J_CDF schemes (Figs. 23 and 24). Figure 25 compares the performance of ASBP and class **D** schemes for parity problems. To the authors' knowledge, the ASBP method is the only method for adaptation of steepness parameter reported in the literature. For all the simulations, we have taken the

momentum factor (α) to be zero. For all simulations, parameters of different schemes are chosen in such a way that best performances will be obtained. The plot for each simulation is averaged over 200 runs.

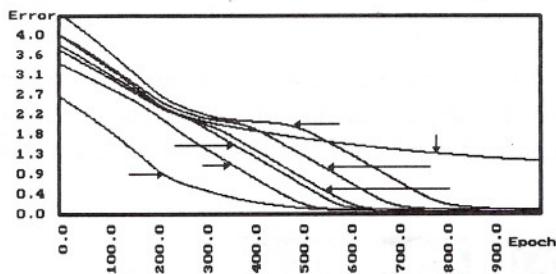


Fig. 11. Digit problem a: standard BP b: Tsetline (4, 4)-AF (η) c: Krinsky(2, 4)-AF (η) d: Krylov (2, 4)-AF (η) e: L_R-p-AV (η) f: VLR g: J(2, 1)-CF (η)

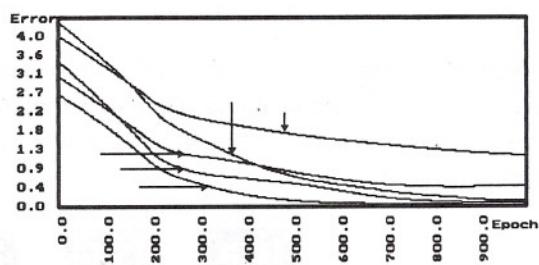


Fig. 12. Digit problem a: standard BP b: Tsetline (4, 4)-CF (η) c: Krinsky(2, 4)-CF (η) d: Krylov (2, 4)-CF (η) e: J(2, 1)-CF (η)

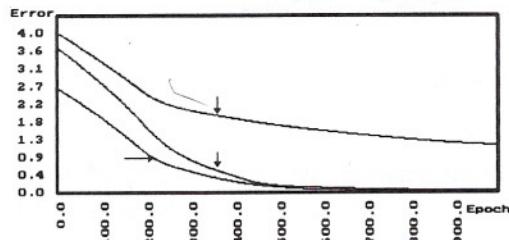


Fig. 13. Digit problem a: Standard BP b: Tsetline(4,6)-Tsetline(2,4)-BF(η) c: J-CF(η) (10, 1)

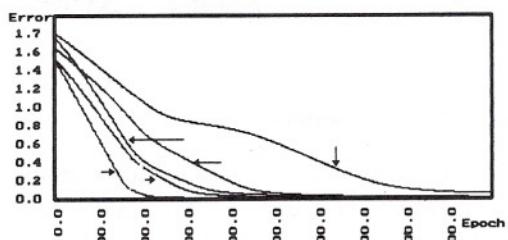


Fig. 14. Parity problem a: Standard BP b: VLR c:J(2, 1)-CF (η) d: J(2, 1)-CF (λ) e: J(2, 1)-J(2, 1)-CDF (η, λ)

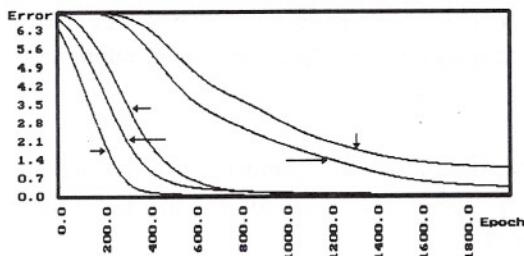


Fig. 15. Encoding problem a: Standard BP b: VLR c:J(2, 1)-CF (η) d: J(5, 6)-CF (λ) e: J(2, 1)-J(5, 6)-CDF (η, λ)

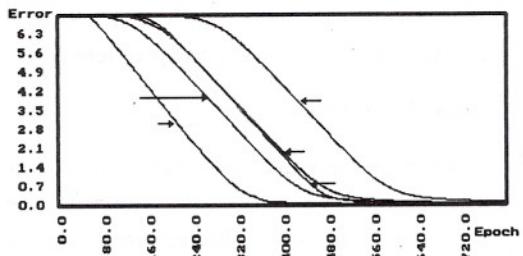


Fig. 16. Symmetry problem a: Standard BP b: VLR c:J(5, 6)-CF (λ) d: J(2, 1)-CF (η) e: J(2, 1)-J(5, 6)-CDF (η, λ)

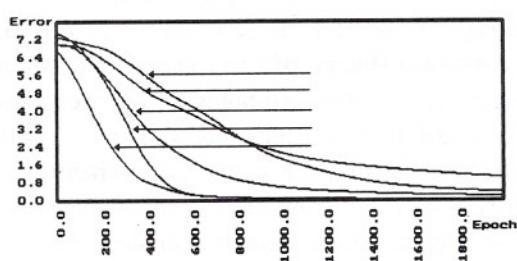


Fig. 17. Encoding Problem a: SC Scheme b: Standard BP c:VLR Scheme d: Fuzzy BP e: J(2, 1)-J(5, 6)-CDF (μ, λ)

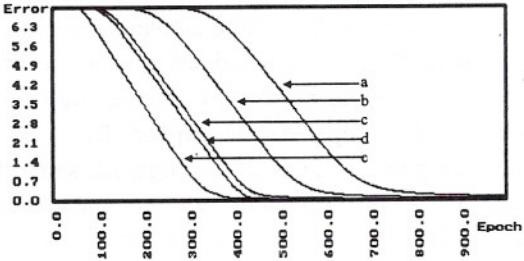


Fig. 18. Symmetry Problem a: SC Scheme b: Standard BP c:VLR Scheme d: Fuzzy BP e: J(2, 1)-J(5, 6)-CDF (μ, λ)

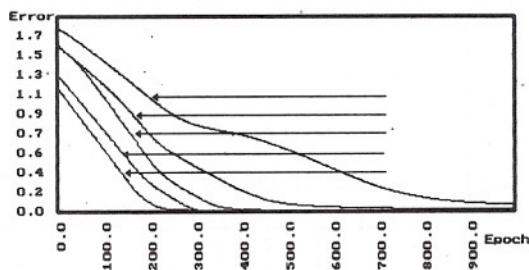


Fig. 19. Parity Problem a: SC Scheme b: Standard BP c: VLR Schemed: Fuzzy BP e: $J(2, 1)$ - $J(2, 1)$ -CDF (μ, λ)

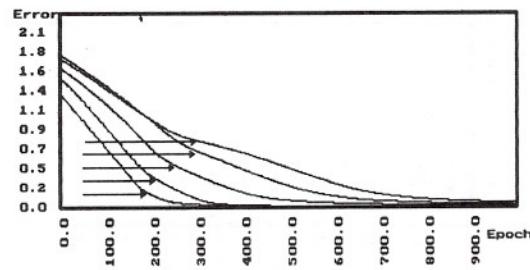


Fig. 20. Parity Problem a:Standard BP b: SAB c: Super SAB d: $J(5, 3)$ -CF (η) e: $J(5, 3)$ - $J(5, 3)$ -CDF(η, λ)

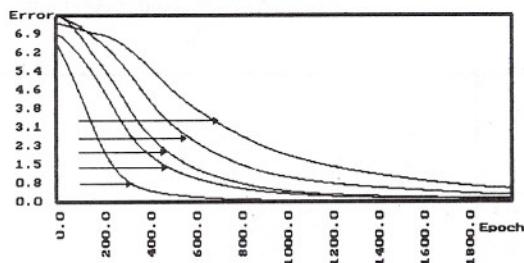


Fig. 21. Encoding Problem a:Standard BP b: SAB c: Super SAB d: $J(5, 3)$ -CF (η) e: $J(5, 3)$ - $J(5, 3)$ -CDF(η, λ)

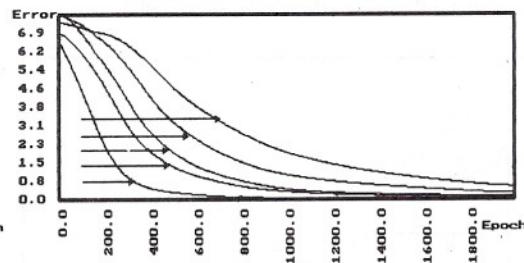


Fig. 22. Digit problem a:Standard BP b: SAB c: Super SAB : $J(5, 3)$ -CF (η)

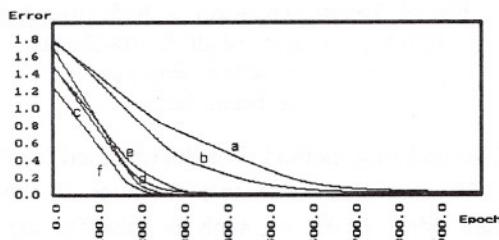


Fig. 23. Parity problem a: Standard BP b: QuickProp c: RPROP d: $J(2,1)$ _DF (λ) e: $J(2,1)$ _CF (η) f: $J(2,1)$ _ $J(2,1)$ -CDF (η, λ)

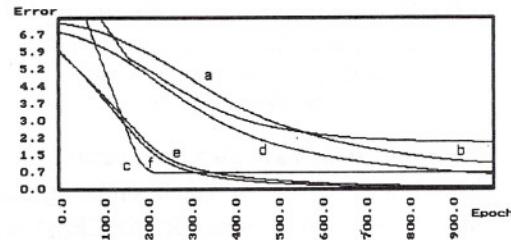


Fig. 24. Encoding problem a: Standard BP b: QuickProp c: RPROP d: $J(5,6)$ _DF (λ) e: $J(2,1)$ _CF (η) f: $J(2,1)$ _ $J(5,6)$ -CDF (η, λ)

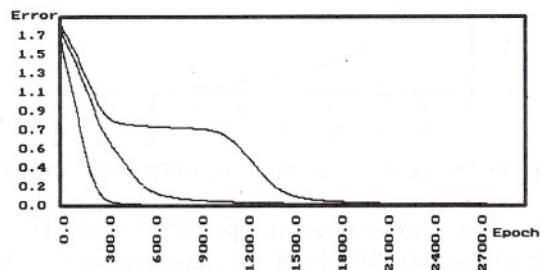


Fig. 25. Parity problem a: Standard BP b: ASBP c: $J(2, 1)$ -DF(λ)

Figures 26 through 29 show the performance of different schemes when both learning rate and momentum factor adapted.

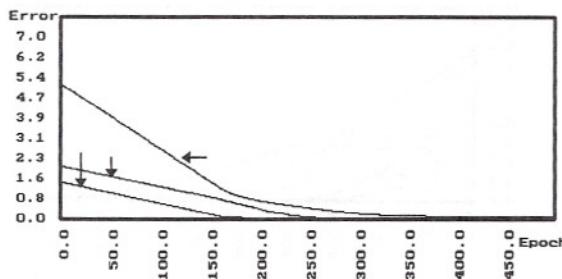


Fig. 26. Digit problem a: $J(5, 10)$ - $CF(\eta)$ b: $J(5, 10)$ - $J(5, 10)$ - $CF(\eta, \alpha)$ c: $J(5, 10)$ - $CF(\eta)$ with constant momentum factor are adapted

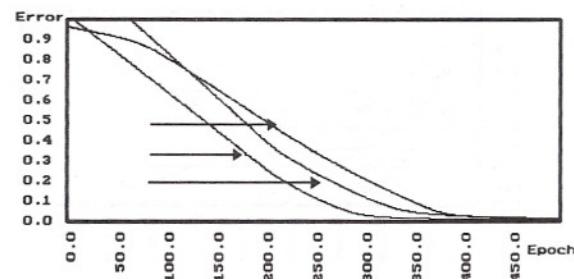


Fig. 27. Parity problem a: $J(5, 10)$ - $J(5, 10)$ - $CDF(\eta, \lambda)$ b: $J(5, 10)$ - $J(5, 10)$ - $CF(\eta, \alpha)$ c: $J(5, 10)$ - $CF(\eta)$ with fixed momentum

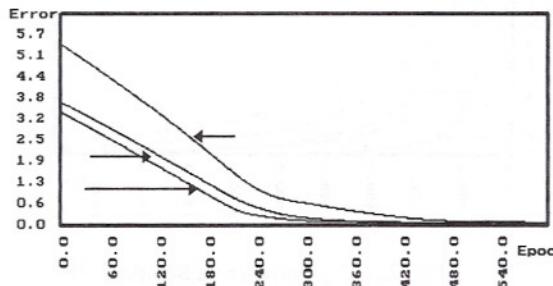


Fig. 28. Encoding problem a: $J(5, 10)$ - $J(5, 6)$ - $CDF(\eta, \lambda)$ b: $J(5, 10)$ - $J(5, 10)$ - $CF(\eta, \alpha)$ c: $J(5, 10)$ - $CF(\eta)$ with fixed momentum factor

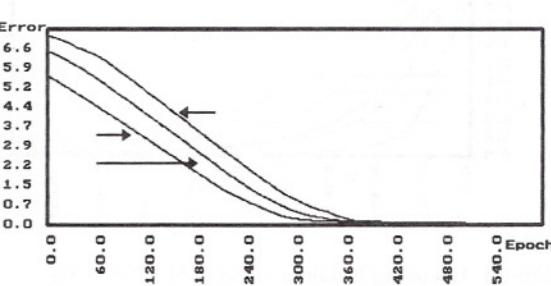


Fig. 29. Symmetry problem a: $J(15, 10)$ - $J(5, 4)$ - $CDF(\eta, \lambda)$ b: $J(15, 10)$ - $J(15, 10)$ - $CF(\eta, \alpha)$ c: $J(15, 10)$ - $CF(\eta)$ with fixed momentum factor

Simulations have also been carried out using on-line updating method. Results obtained for this method are similar to batch updating [26]. Figure 30 compares these two methods for digit recognition problem when $J(5, 6)$ - $CF(\eta)$ scheme is used. More problems, such as classification of sonar signals, vowel recognition, printed Farsi digit recognition, and printed Farsi character recognition problems are tested using LA based schemes which can be found in [27].

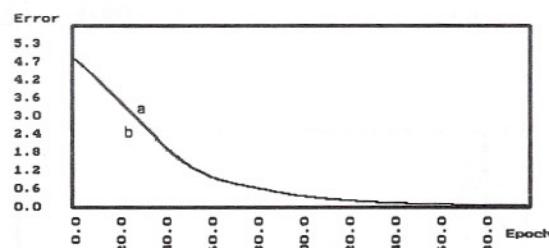


Fig. 30. Digit problem a: on-line updating b: batch updating

Remark 2: J - $CF(\eta)$ scheme has a close relationship with the Jacobs's heuristics. Jacobs [11] has suggested the following heuristics as guidelines for accelerating the convergence of BP learning algorithm through learning rate adaptation.

- Every adjustable network parameter of error function should have its own individual learning-rate parameter.
- Every learning-rate parameter should be allowed to vary from one iteration to the next.

- When the derivative of error function with respect to the synaptic weight has same algebraic sign for several consecutive iterations of the algorithm, the learning-rate parameter for that particular weight should be increased.
- When the algebraic sign of the derivative of error function with respect to the synaptic weight alternates for several consecutive iterations of the algorithm, the learning-rate parameter for that particular weight should be decreased.

Considering the fact that each link of the neural network has its own automata for adaptation of learning parameter, and with the definition of favorable and unfavorable response given before, and also inspecting the transition graphs for these automata, we can see that among different schemes in class CF only J-CF (η) scheme implements all four heuristics of Jacobs. The other schemes in class CF such as Krylov_CF, Krinsky_CF, and Tsetline_CF implement only three of the four heuristics of Jacobs.

Remark 3: J-CF (η) and J-DF (λ) schemes become the standard BP when the memory depth for each action (N) approaches infinity. This is because of the fact that when N is very large it becomes improbable for the $J_{KN,K}$ automata to change action and as a result, a fixed value for the learning rate will be used throughout the training period. Figure 31 shows the effect of memory depth on speed of learning.

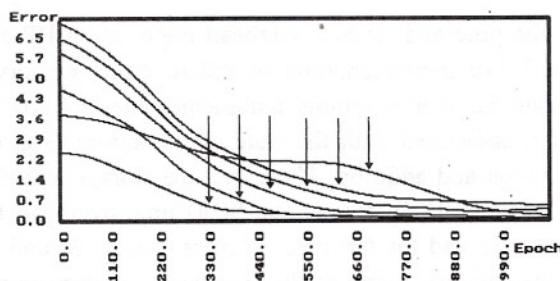


Fig. 31. Digit problem a: Standard BP b: J(10, 10)-CF (η) c: J(10, 250)-CF (η)
d: J(10, 500)-CF (η) e: J(10, 750)-CF (η) f: J(10, 1000)-CF (η)

Remark 4: J-CF (η) and J-DF (λ) schemes become the standard BP when the number of actions (K) approaches infinity. This is due to the fact that when K is very large, the changes in the BP parameters is very small and for a certain amount of changes in the value of the parameter, the automata needs to make a large number of states changes. This effect is the same as the effect we observed for large memory depth and small number of actions. Figure 32 shows the effect of number of actions on speed of learning.

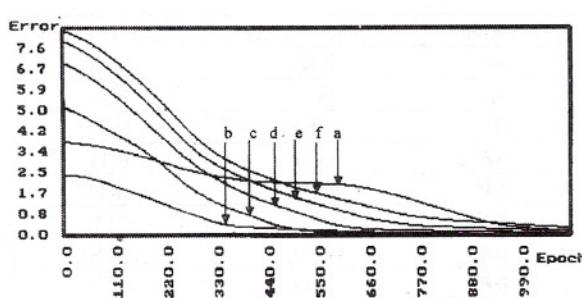


Fig. 32. Digit problem a: Standard BP b: J(10, 10)-CF (η) c: J(250, 10)-CF (η)
d: J(500, 10)-CF (η) e: J(750, 10)-CF (η) f: J(1000, 10)-CF (η)

Remark 5: $J(5,2)$ -DF(λ) ($J(5,2)$ -CF(η)) scheme is used for parity problem, and probability of each action for a randomly selected neuron (weight) is plotted for each scheme. As shown in Figs. 33 and 34 the system converges to the action with the higher value. This effect has also been observed in genetic algorithms for determination of BP parameters [28]. The values of actions are 0.4, 0.8, 1.2, 1.6, and 2.

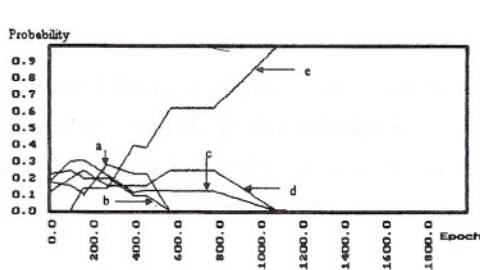


Fig. 33. Action probability of $J(5, 2)$ _DF(λ) scheme for parity problem a: Action 1 b: Action 2 c: Action 3 d: Action 4 e: Action 5

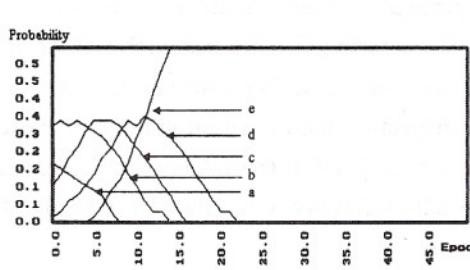


Fig. 34. Action probability of $J(5, 2)$ _CF(η) scheme for parity problem a: Action 1b: Action 2 c: Action 3 d: Action 4 e: Action 5

6. TIME AND SPACE COMPLEXITY OF LA-BASED ADAPTATION SCHEMES

In this section we discuss the time and storage overhead imposed on BP algorithm when LA based adaptation schemes are used. For implementation of FSLA, three memory locations are needed in order to keep track of the state, number of actions, and memory depth of the automata. Changing state and also realizing the action associated with the state of the automata at each epoch, requires few comparisons, integer subtraction and addition. Therefore, the storage and time overhead imposed by each FSLA is $\Theta(1)$. This leads to $\Theta(M)$ storage and $\Theta(M)$ time overhead for **CF** type schemes and $\Theta(N)$ storage and $\Theta(N)$ time overhead for **DF** type schemes, where, M and N are number of weights and number of neurons in the network, respectively. For these schemes, in addition to storage needed to implement FSLA, storage for previous sign of gradient, adapted parameter, and reinforcement signal β for each neuron are needed. When VSLA with K actions used, the storage needed by **CV** and **DV** type schemes are $\Theta(KM)$ and $\Theta(KN)$, respectively. This is because storage is required by each automaton to store its action probability vector P. Due to updating the action probability vector P by the automata at every epoch, the time overhead for **CV** and **DV** type schemes are $\Theta(KM)$ and $\Theta(KN)$, respectively. Class **A** and **B** schemes have overhead of $\Theta(1)$ for both time and storage if FSLA is used, and overhead of $\Theta(K)$ for both time and space if VSLA is used. Table 1 summarizes the storage and time overhead imposed by different schemes.

Table 1: The time and space overhead of proposed schemes

Algorithm	Storage Overhead	Time Overhead
AV	$\Theta(K)$	$\Theta(K)$
AF	$\Theta(1)$	$\Theta(1)$
BV	$\Theta(K)$	$\Theta(K)$
BF	$\Theta(1)$	$\Theta(1)$
CV	$\Theta(KM)$	$\Theta(KM)$
CF	$\Theta(M)$	$\Theta(M)$
DV	$\Theta(KN)$	$\Theta(KN)$
DF	$\Theta(N)$	$\Theta(N)$
SAB	$\Theta(M)$	$\Theta(M)$
SuperSAB	$\Theta(M)$	$\Theta(M)$

In order to justify the overheads imposed by proposed schemes, the ratio of execution time of SAB, SuperSAB, and J-CF schemes to the execution time of standard BP algorithm for parity problem are measured and given in Table 2.

Table 2. Ratio of execution time of different schemes to standard BP

Algorithm	Ratio of Execution Time
SAB	1.66
SuperSAB	1.58
QPROP	1.32
RPROP	1.28
VLR	1.79
CLASS C	1.56
CLASS D	1.55
CLASS CD	1.56

7. CONCLUSION

In this paper a new class of direct methods called **C** class, and a new class of indirect method called **D** class, are presented. All the schemes in these classes use learning automata of fixed or variable structure type to adapt the BP parameter. The adaptation is based on the error surface behavior. To evaluate the performance of these methods, simulation studies were carried out on several learning problems with different error surfaces. Simulations indicate that dynamic adaptation of BP parameters using the proposed methods increase the speed of convergence of the standard BP and is therefore superior to most previous schemes reported for adaptation of BP parameters.

REFERENCES

- Arabshahi, P., Choi, J. J., Marks, R. J. and Caudell, T. P., Fuzzy control of back propagation, Proc. of IEEE Int. Conf. on Fuzzy Systems, p. 967 (1992).
- Kandil, N., Khorasani, K., Patel, R. V. and Sood, V. K., Optimum Learning Rate for Back propagation Neural Networks, In Neural Networks Theory, Technology, and Applications, Edited by P. K. Simpson, p. 249 (1996).
- Parlos, A. G., Fernandez, B., Atya, A. F., Muthusami, J. and Tsai, W. K., An accelerated learning algorithm for multi-layer preceptron networks, *IEEE Trans. on Neural Networks*, **5**, No. 3, p. 493 (1994).
- Cater, J. P., Successfully using peak learning rates of 10 (and greater) in backpropagation networks with the heuristic learning algorithm, *IEEE Proc. of First Int. Conf. on Neural Networks*, **2**, p. 645 (1987).
- Franzini, M. A., Speech Recognition with Backpropagation, IEEE Proc. of Ninth Annual Conf. on Eng. in Medicine and Biology, p. 1702 (1987).
- Vosl, T. P., Mangis, J. K., Rigler, A. K., Zink, W. T. and Alkon, D. L., Accelerating the Convergence of Backpropagation Method, *Biological Cybernetics*, p. 257 (1987).
- Tesauro, G. and Janssens, B., Scaling Relationships in Backpropagation Learning, *Complex Systems*, p. 39 (1988).
- Devos, M. R. and Orban, G. A., Self Learning Backpropagation, Proc. of NeuroNimes (1988).
- Sarkar, D., Methods to speedup error backpropagation learning algorithm, *ACM Computing Surveys*, No. 4, **27**, (1995).
- Tollenaere, T., SuperSAB: fast adaptive backpropagation with good scaling properties, *Neural Networks*, **3** (1990).

11. Jacobs, R. A., Increased rates of convergence through learning rate adaptation, *Neural Networks*, 1, p. 295 (1988).
12. Riedmiller, M. and Heinrich, B., A direct method for faster backpropagation algorithm, *Neural Networks*, 5, p. 465 (1992).
13. Menhaj, M. B. and Hagen, M. H., Rapid learning using modified backpropagation algorithms for multi-layer feedforward neural nets, Proc. of ICEE-95, University of Science and Technology, Tehran, Iran (1995).
14. Sperduti, A. and Starita, A., Speed up learning and network optimization with extended backpropagation, *Neural Networks*, 6, p. 365 (1993).
15. Menhaj, M. B. and Meybodi, M. R., A Novel Learning Scheme for Feedforward Neural Nets, Proc. of ICEE-95, University of Science and Technology, Tehran, Iran (1994).
16. Menhaj, M. B. and Meybodi, M. R., Flexible sigmoidal type functions for neural nets using game of automata, Proc. of Second Annual CSI Computer Conf. CSICC'96, Tehran, Iran, p. 221 (1996).
17. Menhaj, M. B. and Meybodi, M. R., Application of learning automata to neural networks, Proc. of Second Annual CSI Computer Conference CSIC'96, Tehran, Iran, p. 209, (1996).
18. Menhaj, M. B. and Meybodi, M. R., Using learning automata in backpropagation algorithm with momentum, Technical Report, Computer Eng. Department, Amirkabir University of Technology, Tehran, Iran (1997).
19. Beigy, H. and Meybodi, M. R., Adaptation of momentum factor and steepness parameter in backpropagation algorithm using fixed structure learning automata, Proc. of 4th Annual Computer Society of Iran Computer Conference CSICC-99, Sharif University of Technology, Tehran, Iran, p. 117 (1999).
20. Beigy, H., Meybodi, M. R. and Menhaj, M. B., Adaptation of learning rate in backpropagation algorithm using fixed structure learning automata, *Proc. of ICEE-98*, 3, p. 117 (1998).
21. Rumelhart, D. E., Hinton, G. E. and Williams, R. J., Learning Internal Representations by Error Propagation, In parallel distributed processing, Cambridge, MA: MIT Press (1986).
22. Narendra, K. S. and Thathachar, M. A. L., Learning Automata: An Introduction, Prentice-Hall, Englewood cliffs (1989).
23. Meybodi, M. R. and Lakshmivarhan, S., Optimality of a general class of learning algorithm, *Information Science*, 28, p. 1 (1982).
24. Meybodi, M. R. and Lakshmivarhan, S., On a class of learning algorithms which have a symmetric behavior under success and failure, Springer-Verlag, Lecture Notes in Statistics, p. 145 (1984).
25. Meybodi, M. R., Results on strongly absolutely expedient learning automata, Proc. Of OU Inference Conf. 86, ed. D. R. Mootes and Butrick, Athens, Ohio, Ohio Univ. Press, p. 197 (1987).
26. Haykin, S., Neural Networks: A Comprehensive Foundation, Macmillan College Publishing Company, Inc., USA (1994).
27. Meybodi, M. R. and Beigy, H., New classes of learning automata based schemes for adaptation backpropagation algorithm, Technical Reports, Computer Eng. Dept., Amirkabir University of Technology, Tehran Iran (2000).
28. Schaffer, J. D., Whitley, D. and Eshelman, L. J., Combination of genetic algorithms and neural networks: a survey of the state of the art, In Proc. COGANN-92, Int. Workshops on Combination of Genetic Algorithms and Neural Networks, (1992).