

An Analysis of Noise in Recurrent Neural Networks: Convergence and Generalization*

Kam-Chuen Jim

C. Lee Giles*

Bill G. Horne

NEC Research Institute, Inc.
4 Independence Way
Princeton, NJ 08540
{kamjim,giles,horne}@research.nj.nec.com

*Also with
Institute for Advanced Computer Studies
University of Maryland
College Park, MD 20742

Abstract

There has been much interest in applying noise to feedforward neural networks in order to observe their effect on network performance. We extend these results by introducing and analyzing various methods of injecting *synaptic noise* into dynamically-driven recurrent networks during training. We present theoretical results which show that applying a controlled amount of noise during training may improve convergence and generalization performance. In addition, we analyze the effects of various noise parameters (additive vs. multiplicative, cumulative vs. non-cumulative, per time step vs. per string) and predict that best overall performance can be achieved by injecting additive noise at each time step. Noise contributes a second-order gradient term to the error function which can be viewed as an *anticipatory agent* to aid convergence. This term appears to find promising regions of weight space in the beginning stages of training when the training error is large and should improve convergence on error surfaces with local minima. The first order term can be interpreted as a regularization term that can improve generalization. Specifically, this term can encourage internal representations where the state nodes operate in the saturated region of the sigmoid discriminant function. While this effect can improve performance on automata inference problems with binary inputs and target outputs, it is unclear what effect it will have on other types of problems. To substantiate these predictions, we present simulations on learning the dual parity grammar from temporal strings for all noise models, and present simulations on learning a randomly generated 6-state grammar using the predicted best noise model.

Keywords: Recurrent neural networks, synaptic noise, automata, convergence, generalization, regularization, temporal sequences, additive noise, multiplicative noise.

*Printed in *IEEE Transactions on Neural Networks*, vol. 7, no. 6, p.1424-1438, 1996. Copyright IEEE.

1 INTRODUCTION

There has been a great deal of work on enhancing neural network performance. Two important performance parameters are convergence and generalization. Convergence is the amount of time, measured in CPU operations or training epochs, required to find an acceptable solution during training. Generalization measures the ability to correctly classify new unseen data.

Among some of the important methods that have been introduced to improve convergence is the notion of searching the coarse regions of state space before searching finer regions. *Simulated annealing* is one such method, and was shown by Kirkpatrick, et.al. [20] to be very effective on combinatorial problems like graph partitioning. Similarly, Ackley, et.al. and Hinton and Sejnowski [1, 14] used noise in Boltzmann machines to escape from local minima. Specifically, it was indicated that at high noise temperatures Boltzmann machines will search the overall structure of state space for a coarse minimum, while subsequently lowering the noise temperature will allow a better minimum to be found within the coarse-scale minimum.

It has been shown that generalization can be improved by removing the redundancies and irrelevancies in a network. There have been at least two approaches: pruning and regularization. An example of pruning is *optimal brain damage*, which can improve generalization ability and speed of learning by using second-derivative information to remove unimportant weights from the network (Cun, Denker, and Solla [8]). Also, Giles and Omlin [12] have demonstrated improvement in generalization of recurrent neural networks after pruning. Other pruning methods can be found in Reed [26]. *Weight decay* is an example of a regularization method, and was shown by Krogh and Hertz [21] to improve generalization on feed-forward networks by suppressing irrelevant components of the weight vector.

1.1 Previous Work on Training with Noise

Previous research has investigated the effects of *noise* on feedforward neural networks. Training with noise can lead to more realistic biological models. For example, Hinton, Sejnowski, and Ackley [15] compared noise temperature in the Boltzmann machine to the addition of Gaussian noise at neuronal membrane potentials. Bremermann and Anderson [4, 5] introduced a biased random walk in weight space (the "chemotaxis algorithm") as a biologically plausible learning rule, and Kilis and Ackerman [19] introduced a *stochastic neuron* model which uses Monte Carlo techniques to stochastically bias the firing decision of neurons. More relevant to this paper however, are the experiments that use noise injection during training to enhance network performance.

Much work in training with noise has been motivated by the desire to improve *fault tolerance*, i.e. the

graceful degradation of performance in the presence of faults. Judd and Munro [18] showed that using noisy hidden nodes during training can result in error-correcting codes which increase the tolerance of feedforward nets to unreliable nodes. Similarly, Séquin and Clay [28] showed that randomly disabling hidden nodes (i.e. clamp to zero) during the training phase increases the tolerance of multilayer perceptrons to node failures. Minnix [23] also demonstrated improved fault tolerance by training on noisy inputs, and argued that such a method helps prevent *pattern memorization*. Murray and Edwards [24] utilized *synaptic noise* during training to improve fault tolerance and training quality. Omlin and Giles [25] have extended the fault tolerance considerations to recurrent neural networks.

With the exception of Murray and Edwards [24], these methods generally indicate convergence as a trade-off for enhanced performance. Recently, stochastic techniques, such as Nodal Perturbation, Weight Perturbation, and Summed Weight Neuron Perturbation, have been introduced which in some cases may converge more quickly than pure gradient descent (Cauwenberghs [7], Dembo and Kailath [9], Flower and Jabri [10, 17]). These methods in general require more epochs to converge but less computation per epoch, mainly because calculation of the true gradient is not required. Also, it has been demonstrated that faster learning rates can be achieved by adding noise to the weight updates and adapting the magnitude of such noise to the output error (Burton and Mpitsos [6]). Hanson developed a stochastic version of the delta rule which adapt weight means and standard deviations instead of clean weight values, and demonstrated improved convergence rates [13].

Training with noise can also improve generalization. Holmström and Koistinen [16] showed that training with additive input noise can improve generalization in feedforward networks, and provided mathematically justified rules for choosing the characteristics of the noise. Bishop [3] showed that training with noisy data is equivalent to Tikhonov Regularization and suggested directly minimizing the regularized error function as a practical alternative.

1.2 Overview of paper

To date, there has been little work on applying noise insertion methods to recurrent neural networks. In this paper, we explore the effects of synaptic noise on recurrent neural networks. We introduce several ways to inject noise during training, and demonstrate that utilizing *noise-enhanced* error functions can improve convergence and generalization. Previous work on improving these two performance measures, as discussed above, focused on simplifying feedforward neural networks and on ways of first searching the coarse regions of state space. The function of synaptic noise in recurrent nets is novel in the sense that synaptic noise can improve convergence by searching for *promising regions* of state space, and at the same

time enhance generalization by favoring *saturated* state nodes. Promising regions of state space do not necessarily correspond to coarsely-spaced local minima, and saturated state nodes are in general different from a simplified network.

Besides an analysis of noise in recurrent neural networks, we perform many simulations to verify our analysis. For recurrent neural network and problem selection, we use a second-order fully-connected recurrent neural network. This recurrent neural network is then trained to learn automata from grammatical strings encoded as temporal sequences. For a small neural network and small automata, we ran many simulations and thus show statistical confidence in our results.

The next section, section 2, defines the second-order fully-connected recurrent neural network used in simulations and introduces several methods of implementing synaptic noise. We discuss different noise models and compare their effects on the error function in section 3. From this analysis, predictions are made regarding their effects on convergence and generalization performance. Section 4 presents simulation results on learning the dual parity and a small 6-state randomly generated automata from strings encoded as temporal sequences. These results are general to all incremental weight updating methods and should extend easily to other recurrent networks. Conclusions and open questions follow in section 5.

2 NOISE INJECTION IN RECURRENT NEURAL NETWORKS

2.1 Noiseless Recurrent Network Architecture and Training

In this paper we restrict ourselves to second-order, fully-connected, discrete-time recurrent networks since such networks have demonstrated good performance on the types of problems we discuss in this paper, namely, that of learning finite state grammars [11, 22, 29]. The discrete-time recurrent network architecture consists of N recurrent state neurons S_j and L non-recurrent input neurons I_k . The states are related by a set of N^2L weights W_{ijk} . The recurrent network operation resembles that of a finite state machine with the state process: {input, current state} \implies {next state}, as defined by the following equations:

$$S_i^{t+1} = g(\Theta_i^{t+1}) \tag{1}$$

$$\Theta_i^{t+1} = \sum_{j,k} W_{ijk} S_j^t I_k^t \tag{2}$$

where Θ_i is the net input to unit i , and $g(x)$ is a predetermined sigmoid discriminant function. In the simulations described in section 4, $g(x)$ is defined as

$$g(x) = \frac{1}{1 + e^{-x}} \quad (3)$$

“One-hot” or “unary” input encoding is used. The input symbols are encoded as orthonormal unit vectors in L -dimensional space, so the inputs are in the form $I_k = \delta_{k\alpha}$ for $(0 \leq \alpha \leq L)$. For the dual parity DFA, $L = 3$, in order to represent the binary input characters 0, 1, and the empty string **e**. Hence, the possible symbols are (0,0,1), (0,1,0), and (1,0,0).

Each input string is presented to the network one character per discrete time step t . At time $t = T$, when the entire input string has been presented to the network, the value at a selected neuron S_O is defined to be the output. An input string is accepted if $S_O > 1 - \epsilon$, or rejected if $S_O < \epsilon$, where ϵ is the *error tolerance* and is a positive real number. Weights are updated at the end of each string.

During training, an error function is computed as

$$E_p = \frac{1}{2}\epsilon_p^2 = \frac{1}{2}(S_O^T - d_p)^2 \quad (4)$$

where d_p is the target output value for pattern p (either 1 or 0 depending on the DFA response of the input string), and ϵ_p is the raw error.

For complexity reasons we used the Back-Propagation-Through-Time (BPTT) training algorithm [27]. Incremental learning, as described in Giles, et.al. [11], is used to descend the error function. Basically, the network is trained on a subset of the training set, called the *working set*, which gradually increases in size until the network is able to correctly classify the entire training set. Strings from the working set are presented to the network in *alphabetical order*. Incremental learning reduces the training time required because often not all the strings in the training set need to be seen before the network can classify the entire set correctly. Training stops when all strings in the training set are correctly classified.

2.2 Noise Injection

Murray and Edwards [24] simulated synaptic noise on multilayered perceptrons (MLPs) by applying noise to the weights of each layer during training. Applying this method to recurrent networks is not straightforward because effectively the same weights are propagated forward in time. This can be seen by recalling the BPTT representation of *unrolling* a recurrent network in time into T layers with identical weights, where T is the

length of the temporal sequence. As in Burton and Mpitsos [6], Hanson [13], and Murray and Edwards [24] we will only focus on noise injection in synaptic weights and ignore other methods of noise insertion.

Because of its temporal nature, synaptic noise injection of recurrent networks introduces two issues which do not pertain to feedforward nets:

1. What is the frequency of noise injection?
2. Is the noise injection cumulative or non-cumulative?

This is a direct result of the fact that recurrent networks effectively propagate the same weights forward in time. Should one set of noise values also be propagated forward, or should a different set of noise values be applied at each time step? Or, should noise be applied only at the beginning of each input sequence? If noise is propagated forward in time, should it accumulate in the weights over all time steps, or can the noise be inserted independently each time step?

In this paper we consider the cases where noise is propagated forward in time. We introduce a *per time step* method of injecting synaptic noise into recurrent networks which inserts a new random noise value at each time step of the forward pass during training. On the other hand, a *per string* model can be viewed as a special case of the per time step model such that the same set of noise values is injected for all time steps of each given input string.

Non-cumulative models apply a new noise value Δ_{ijk} to a noise-free weight at each instance of noise injection, while cumulative models apply a new noise to a previously perturbed weight at each instance. Non-cumulative noise can be represented as

$$W_{t,ijk} = N(W_{ijk}^*), \quad (5)$$

where $\{W_{ijk}^*\}$ is the noise-free weight set and N is the noise process applied to weight W_{ijk}^* . Cumulative noise be represented by

$$W_{0,ijk} = N(W_{ijk}^*), \quad (6)$$

$$W_{t,ijk} = N(W_{t-1,ijk}), t > 0. \quad (7)$$

An issue common to both feedforward and recurrent networks is whether noise is additive or multiplicative. Additive noise perturbs a weight parameter such that

$$W_{ijk} \longrightarrow W_{ijk} + \Delta_{ijk}, \quad (8)$$

	t_1	t_2	t_3	\dots
per time step non-cumulative	$W^* + \Delta_1$	$W^* + \Delta_2$	$W^* + \Delta_3$	\dots
per time step cumulative	$W^* + \Delta_1$	$W^* + \Delta_1 + \Delta_2$	$W^* + \Delta_1 + \Delta_2 + \Delta_3$	\dots
per string non-cumulative	$W^* + \Delta_1$	$W^* + \Delta_1$	$W^* + \Delta_1$	\dots
per string cumulative	$W^* + \Delta_1$	$W^* + 2\Delta_1$	$W^* + 3\Delta_1$	\dots

Table 1: Sample noise injection steps for all additive noise models at time steps t_1 , t_2 , and t_3 . W^* is the noise-free weight. Δ 's are the noise terms for each time step.

where Δ_{ijk} is a random noise value. On the other hand, multiplicative noise perturbs each weight such that

$$W_{ijk} \longrightarrow W_{ijk} + W_{ijk}\Delta_{ijk}. \quad (9)$$

Table 1 clarifies these issues by illustrating the noise injection steps for all additive noise processes. The Δ 's are the various noise terms, and W^* is the noise free-weight. For per time step non-cumulative noise insertion a different noise is added at each time step and the noise from the previous time step does not accumulate in time. For the per time step cumulative noise insertion all previously added noise is carried over to the next time step. For per string (or per sequence) non-cumulative, the same noise is inserted at the beginning and used throughout the temporal updates. For per string cumulative, the initial noise accumulates for each time step. Similar results would be seen for the multiplicative noise case.

We report experiments which explore all combinations of additive vs. multiplicative, cumulative vs. non-cumulative, and per string vs. per time step noise models. In all noise models the actual weight updates ΔW_{ijk} 's descend the noise-enhanced error term E_p . These updates are applied to the noise-free weight set. This prevents "noisy updates", allows simpler mathematical analysis, and for the purpose of these experiments has the desirable effect of limiting the effect of noise to the error function.

For all cases the noise source is chosen to be a random, zero-mean white noise process uniformly distributed between $(-\mathbf{a}, \mathbf{a})$, where \mathbf{a} is a positive constant parameter.

3 ANALYSIS OF SYNAPTIC NOISE

The effects of each noise model is analyzed by taking the Taylor expansion on the error function around the noise-free weight set. By truncating this expansion to second order and lower terms, we can interpret the effect of noise as a set of regularization terms applied to the error function. From these terms predictions can be made about the effects on generalization and convergence. A similar analysis was performed by Murray and Edwards [24] on MLPs to demonstrate the effects of synaptic noise on fault tolerance and training quality. Le Cun et. al. [8] also applied a Taylor expansion on the error function to approximate the effects

	NON-CUMULATIVE	
	Additive	Multiplicative
Noise step	$W_{t,ijk} = W_{ijk}^* + \Delta_{t,ijk}$	$W_{t,ijk} = W_{ijk}^* + W_{ijk}^* \Delta_{t,ijk}$
1st order	$\frac{1}{2}\sigma^2 \sum_{t=0}^{T-1} \sum_{ijk} \left(\frac{\partial S_O^T}{\partial W_{t,ijk}} \right)^2$	$\frac{1}{2}\sigma^2 \sum_{t=0}^{T-1} \sum_{ijk} \left(W_{t,ijk} \frac{\partial S_O^T}{\partial W_{t,ijk}} \right)^2$
2nd order	$\frac{1}{2}\epsilon_p \sigma^2 \sum_{t=0}^{T-1} \sum_{ijk} \frac{\partial^2 S_O^T}{\partial (W_{t,ijk})^2}$	$\frac{1}{2}\epsilon_p \sigma^2 \sum_{t=0}^{T-1} \sum_{ijk} (W_{t,ijk})^2 \frac{\partial^2 S_O^T}{\partial (W_{t,ijk})^2}$
	CUMULATIVE	
	Additive	Multiplicative
Noise step	$W_{t,ijk} = W_{ijk}^* + \sum_{\tau=0}^t \Delta_{\tau,ijk}$	$W_{t,ijk} = W_{ijk}^* \prod_{\tau=0}^t (1 + \Delta_{\tau,ijk})$
1st order	$\frac{1}{2}\sigma^2 \sum_{t,u=0}^{T-1} \sum_{ijk} v \frac{\partial S_O^T}{\partial W_{t,ijk}} \frac{\partial S_O^T}{\partial W_{u,ijk}}$	$\frac{1}{2}\sigma^2 \sum_{t,u=0}^{T-1} \sum_{ijk} v W_{t,ijk} W_{u,ijk} \frac{\partial S_O^T}{\partial W_{t,ijk}} \frac{\partial S_O^T}{\partial W_{u,ijk}}$
2nd order	$\frac{1}{2}\epsilon_p \sigma^2 \sum_{t,u=0}^{T-1} \sum_{ijk} v \frac{\partial^2 S_O^T}{\partial W_{t,ijk} \partial W_{u,ijk}}$	$\frac{1}{2}\epsilon_p \sigma^2 \sum_{t,u=0}^{T-1} \sum_{ijk} v W_{t,ijk} W_{u,ijk} \frac{\partial^2 S_O^T}{\partial W_{t,ijk} \partial W_{u,ijk}}$

Table 2: *Per time step* noise injection terms and error function expansion terms for cumulative and non-cumulative noise for both the additive and multiplicative cases. $v = \min(t+1, u+1)$. W^* is the noise-free weight set.

of small weight perturbations, and dropped the first order terms by assuming the network is at convergence. As in [24], we do not assume the networks are at convergence and consider both first and second order terms. Tables 2 and 3 list the noise injection step and resulting first and second order Taylor expansion terms for all noise models. These terms are derived in Appendix A.

Our results are comparable (in terms of the effects on the error function) to feedforward nets only when the recurrent noise models are non-cumulative and per time-step. Cumulative noise and per string noise in recurrent nets are not comparable to their feedforward versions because of the absence of error-weight correlations across time in feedforward nets. In feedforward nets, per string noise is identical to per time step noise because the length of each pattern in time is always one. Also, cumulative noise in feedforward networks are identical to noncumulative noise when the weights are updated after the presentation of each training sample, but have different variance when using batch training (see Appendix B). Thus, our non-cumulative, multiplicative, per time step analytical results are consistent with Murray and Edwards [24], who applied multiplicative synaptic noise to MLPs. Other noise models accumulate noise at different rates and have different convergence and generalization effects, as described below and summarized in Tables 2, 3, and 4.

	NON-CUMULATIVE	
	Additive	Multiplicative
Noise step	$W_{t,ijk} = W_{ijk}^* + \Delta_{ijk}$	$W_{t,ijk} = W_{ijk}^* + W_{ijk}^* \Delta_{ijk}$
1st order	$\frac{1}{2}\sigma^2 \sum_{t,u=0}^{T-1} \sum_{ijk} \frac{\partial S_O^T}{\partial W_{t,ijk}} \frac{\partial S_O^T}{\partial W_{u,ijk}}$	$\frac{1}{2}\sigma^2 \sum_{t,u=0}^{T-1} \sum_{ijk} W_{t,ijk} W_{u,ijk} \frac{\partial S_O^T}{\partial W_{t,ijk}} \frac{\partial S_O^T}{\partial W_{u,ijk}}$
2nd order	$\frac{1}{2}\epsilon_p \sigma^2 \sum_{t,u=0}^{T-1} \sum_{ijk} \frac{\partial^2 S_O^T}{\partial W_{t,ijk} \partial W_{u,ijk}}$	$\frac{1}{2}\epsilon_p \sigma^2 \sum_{t,u=0}^{T-1} \sum_{ijk} W_{t,ijk} W_{u,ijk} \frac{\partial^2 S_O^T}{\partial W_{t,ijk} \partial W_{u,ijk}}$
	CUMULATIVE	
	Additive	Multiplicative
Noise step	$W_{t,ijk} = W_{ijk}^* + (t+1)\Delta_{ijk}$	$W_{t,ijk} = W_{ijk}^* (1 + \Delta_{ijk})^t$
1st order	$\frac{1}{2}\sigma^2 \sum_{t,u=0}^{T-1} \sum_{ijk} (t+1)(u+1) \frac{\partial S_O^T}{\partial W_{t,ijk}} \frac{\partial S_O^T}{\partial W_{u,ijk}}$	$\frac{1}{2}\sigma^2 \sum_{t,u=0}^{T-1} \sum_{ijk} (t+1)(u+1) \frac{\partial S_O^T}{\partial W_{t,ijk}} \frac{\partial S_O^T}{\partial W_{u,ijk}}$ $+ \epsilon_p \sigma^2 \sum_{t=0}^{T-1} \sum_{ijk} t(t+1) \frac{\partial S_O^T}{\partial W_{t,ijk}}$
2nd order	$\frac{1}{2}\epsilon_p \sigma^2 \sum_{t,u=0}^{T-1} \sum_{ijk} (t+1)(u+1) \frac{\partial^2 S_O^T}{\partial W_{t,ijk} \partial W_{u,ijk}}$	$\frac{1}{2}\epsilon_p \sigma^2 \sum_{t,u=0}^{T-1} \sum_{ijk} (t+1)(u+1) \frac{\partial^2 S_O^T}{\partial W_{t,ijk} \partial W_{u,ijk}}$

Table 3: *Per string* noise injection terms and error function expansion terms for cumulative and non-cumulative noise for both the additive and multiplicative cases. W^* is the noise-free weight set.

3.1 Generalization Improvement due to Synaptic Noise

We propose that two ways to improve generalization are

- Reduce the raw error, and
- Improve the internal representation of information in the network.

Reducing the raw error ϵ_p improves the tolerance of the network to perturbations in the output nodes. This tolerance provides a buffer from which slight variations in the output node will lead to the same classification, thus improving generalization. At the same time, it has been shown that on networks with too many parameters over-training can lead to poor generalization, since reducing the raw error to zero causes the network to effectively memorize the training set [2].

One common cause of poor generalization in recurrent networks that are trained on hidden-state problems is the presence of unsaturated state representations [30]. Typically, a network cannot revisit the exact same point in state space, but tends to wander away from its learned state representation. One approach to alleviate this problem is to encourage state nodes to operate in the saturated regions of the sigmoid. The first order error expansion terms of most noise models considered are capable of encouraging the network to achieve saturated states. This can be shown by applying the chain rule to the partial derivative in the first

order expansion terms:

$$\frac{\partial S_O^T}{\partial W_{t,ijk}} = \frac{\partial S_O^T}{\partial \Theta_O^T} \sum_l \left[\left(\frac{\partial \Theta_O^T}{\partial S_l^{T-1}} \frac{\partial S_l^{T-1}}{\partial \Theta_l^{T-1}} \right) \cdots \sum_n \left(\frac{\partial \Theta_m^{t+1}}{\partial S_n^t} \frac{\partial S_n^t}{\partial W_{t,ijk}} \right) \right], \quad (10)$$

where Θ_i^t is the net input to state node i at time step t . The partial derivatives $\frac{\partial S}{\partial \Theta}$ favor internal representations such that the effects of perturbations at the net inputs Θ_i^t on the states S_i^t are minimized. While this effect may improve performance on automata inference problems (learning automata from temporal sequences of strings) with binary inputs and target outputs, it is unclear what effect it will have on other types of problems.

Multiplicative noise implements a form of weight decay because the error expansion terms include the weight products $W_{t,ijk}^2$ or $W_{t,ijk}W_{u,ijk}$. Although weight decay has been shown to improve generalization on feedforward networks (Krogh and Hertz, 1992) we hypothesize that in general weight decay will not always improve generalization for recurrent networks. that are learning finite state automata (FSA) problems. Large weights are necessary to saturate the state nodes to the upper and lower limits of the sigmoid discriminant function. Therefore, we predict additive noise will allow better generalization on automata inference problems because it is not a weight decay term, while multiplicative noise will work better on other types of problems where weight decay is appropriate.

Noise models whose first order error term contain the expression $\frac{\partial S_O^T}{\partial W_{t,ijk}} \frac{\partial S_O^T}{\partial W_{u,lmn}}$ will favor saturated states for those partials whose sign correspond to the sign of a majority of all partials. These noise models will favor unsaturated states, operating in the linear region of the sigmoid, for partials whose sign is in the minority. Such *sign-dependence* is not optimal because not all the weights are encouraged to be saturated.

The error terms for cumulative per time step noises sum a product with the expression $v \frac{\partial S_O^T}{\partial W_{t,ijk}} \frac{\partial S_O^T}{\partial W_{u,lmn}}$, where $v = \min(t+1, u+1)$. The effect of cumulative noise increases more rapidly because of v and thus optimal generalization and detrimental noise effects will occur at lower amplitudes than non-cumulative noise.

For cumulative per string noise models, the products $(t+1)(u+1)$ and $\Psi_{t,ijk}\Psi_{u,lmn}$ in the expansion terms rapidly overwhelm the raw error term. Generalization improvement is not expected for these models.

We also reason that all generalization enhancements will be valid only for a range of noise values, above which noise overwhelms the raw error information. The differences between the noise models are summarized in Table 4, where a check indicates that the corresponding enhancement is predicted. It is important to realize that although all noise models are theoretically capable of improving convergence and generalization by searching for promising weights and saturated states, Table 4 lists the *additional* enhancements to the

	Generalization Improvement	Convergence Improvement	Generalization & Convergence Improvement
	Sign-independent favoring of saturated states	Anticipatory effect	No additional weight terms
Per Timestep Cumulative			
Additive	✓	✓	✓
Multiplicative			
Per Timestep Non-cumulative			
Additive	✓	✓	✓
Multiplicative	✓	✓	
Per String Cumulative			
Additive		✓	✓
Multiplicative			
Per String Non-cumulative			
Additive		✓	✓
Multiplicative			

Table 4: Generalization and Convergence Improvement: Checks indicate additional enhancements to the search for saturated states/promising weights due to each noise model. Sign-independent enforcement of saturated states improves generalization (column 2). The anticipatory effect is guaranteed only when there are either no weight terms in the second order derivatives of Tables 2 and 3, or the weight terms are always positive (column 3). The absence of weight terms in Tables 2 and 3 improves convergence/generalization by not constraining the weights to small amplitudes (column 4).

search. From this table it can be inferred that additive per time step noise models should yield the best generalization performance.

3.2 Convergence Improvement due to Synaptic Noise

Convergence can be represented by the number of epochs required to learn the training set. This measure is dependent on the error trajectory followed during the training phase. On error surfaces filled with local minima many epochs may be required before the training set is learned, because the learning trajectory may enter an unacceptable local minima and face a difficult or impossible task of exiting from it.

Intuitively, synaptic noise may alleviate this problem by allowing the network to make noisy jumps out of local minima. Obviously, this will be beneficial only to learning problems whose error surfaces contain local minima. On simple error surfaces synaptic noise may worsen generalization performance by perturbing a straightforward training trajectory which would have otherwise led directly to a solution.

Synaptic noise can improve convergence by favoring *promising* weights in the beginning stages of training. This can be demonstrated by examining the second order error expansion term for non-cumulative, multiplicative, per time step noise:

$$\frac{1}{2} \epsilon_p \sigma^2 \sum_{t=0}^{T-1} \sum_{ijk} (W_{t,ijk})^2 \left(\frac{\partial^2 S_O^T}{\partial (W_{t,ijk})^2} \right).$$

When ϵ_p is negative, solutions with a negative second order state-weight partial derivative will be destabilized. In other words, when the output S_O^T is too small the network will favor updating in a direction such that the first order partial derivative is increasing. A corresponding relationship can be observed for the case when ϵ_p is positive. Thus the second order term of the error function will allow a higher raw error ϵ_p to be favored if such an update will place the weights in a more promising area, i.e. a region where weight changes are likely to move S_O^T in a direction to reduce the raw error. The *anticipatory effect* of this term, or *look-ahead* property as described in Murray and Edwards [24], is more important in the beginning stages of training where ϵ_p is large, and will become insignificant in the finishing stages of training as ϵ_p approaches zero.

In general, the second order terms of all noise models, as shown in Tables 2 and 3, are capable of this anticipatory effect. In all noise models, these terms will favor updating the weights such that the sum of the dependencies of the output node on all weights has the opposite sign of the raw error term ϵ_p . In other words, either a *majority* of the weights are moved to a more *promising* space, or a minority of the weights have moved to a more promising space with very steep slopes. Non-cumulative, multiplicative, per time step noise includes the squared terms $W_{t,ijk}^2$, which are always positive. Thus, this noise model and all additive noise models are guaranteed to have this anticipatory effect. All other multiplicative noise models have the products $W_{t,ijk}W_{u,ijk}$ or $W_{t,ijk}W_{u,lmn}$, implying that the anticipatory effect is observed only when the majority of the weights are positive. This is summarized in Table 4, column 3.

Similar to arguments in Section 3.1, the absence of weight decay will make the learning task easier and improve convergence (see Table 4, column 4).

4 SIMULATION RESULTS

4.1 Experimental Methodology

Two sets of experiments which explore the effects of all combinations of cumulative vs. non-cumulative, additive vs. multiplicative, and per string vs. per time step noise are reported in this paper. In order to perform many experiments in a reasonable amount of computation time, we attempt to learn the simple dual parity grammar from sample strings encoded as temporal sequences. Dual parity is a 4-state automaton that recognizes binary strings containing an even number of ones and zeroes (see figure 1) and is a classic example of a “hidden-state” problem. Both sets of experiments were performed on the dual parity problem, one set using recurrent networks with 3 state neurons, and the other using recurrent networks with 4 state neurons. In addition, another set of simulations apply the best predicted noise model (additive, non-cumulative, per

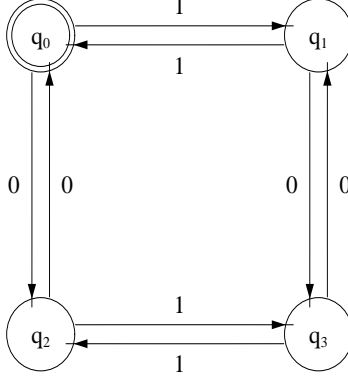


Figure 1: The Dual Parity Finite State Automaton. State q_0 is both start and final state.

time step) to learning a slightly larger, randomly generated 6 state automaton (RND6 – see figure 2). Our experiments consist of 500 simulations for each data point and achieve useful (90%) confidence levels.

Second order recurrent networks with 3 input neurons are trained on the dual parity problem. The finite state automaton (FSA) that recognizes the dual parity grammar is shown in figure 1. The data consists of 8191 strings of lengths 0 to 12. The training set consists of the first 1023 strings of lengths 0 to 9, while the initial working set consists of 31 strings of lengths 0 to 4. Generalization performance of a trained network was obtained by testing on the entire data set of 8191 strings. Training was stopped when all of the training set was correctly classified to the previously discussed error criterion. If training persisted beyond 5000 epochs, training was halted and the trained network was not used. Thus, only networks which were perfectly trained on the training set were tested for generalization. For runs where the noise amplitude was small, the networks usually converged in less than 5000 epochs. For higher noise values, up to a quarter of the runs did not correctly classify the training set in the required number of epochs and were discarded.

Second order recurrent networks with 3 input neurons and 10 state neurons were trained on the RND6 grammar. The finite state automaton that recognizes this grammar is shown in figure 2. The data consists of 1023 strings of length 0 to 9. The training set consists of the first 50 strings. If training persisted beyond 500 epochs, training was halted and the trained network was not used.

As mentioned before the noise source is chosen to be a random, zero-mean white noise process uniformly distributed between $(-\mathbf{a}, \mathbf{a})$, where \mathbf{a} is the amplitude.

To ensure consistency and fairness across all experiments, the learning rate and momentum were set to 0.5, and the weights were initialized to random values between $[-1.0, 1.0]$ in all simulations. The empirical results presented below represent 90% confidence levels using a 2-tailed Student's t -distribution, and were obtained by running 500 simulations for each data point.

The convergence and generalization results for the dual parity problem are presented first, followed by

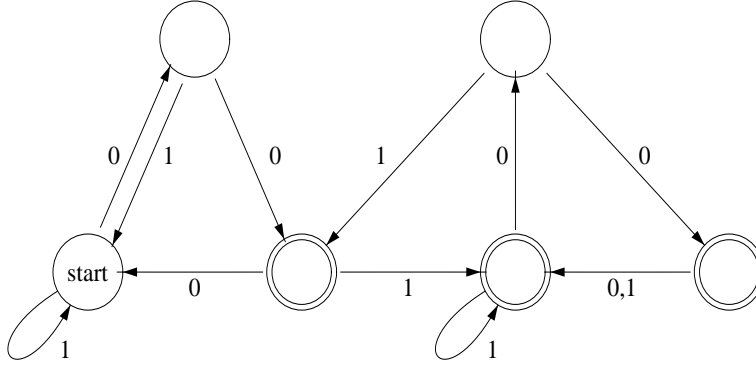


Figure 2: The Randomly Generated 6 State Automaton. Double circles are the accepting states.

the results for RND6.

4.2 Generalization Performance

In figures 3 and 4, we see generalization performance on the dual parity problem as a function of noise amplitude for both the 3-state and 4-state neuron recurrent nets. For each net we see generalization for various noise combinations. For all simulations the zero-noise case is where $a = 0$. Generalization performance mirrors some of our predictions. However, all curves exhibit the same general shape of generalization performance. The generalization improves as noise increases from the zero-noise point for all noise cases and then starts to degrade after a certain noise-threshold amplitude. After this noise-threshold amplitude, generalization performance eventually becomes worse than the zero-noise case. (The curve in Figure 3a does increase as a increases, but is not shown.) The cumulative per string noise cases for both additive and multiplicative noise failed to converge for all runs.

4.3 Convergence Time

Convergence time on the dual parity problem as a function of noise amplitude are plotted in figures 5 and 6. As mentioned above, networks injected with cumulative per string noise do not converge. For most noise models, convergence for added noise initially decreases up to a noise-threshold and then starts to increase and finally exceeds that for the no-noise case. Similar behavior for feedforward networks has been observed [24]. This is more apparent for the 3-state neuron simulations than the 4-state neuron ones; however, close inspection will show there is a convergence mean less than that for the zero-noise case for the 4-state neuron case.

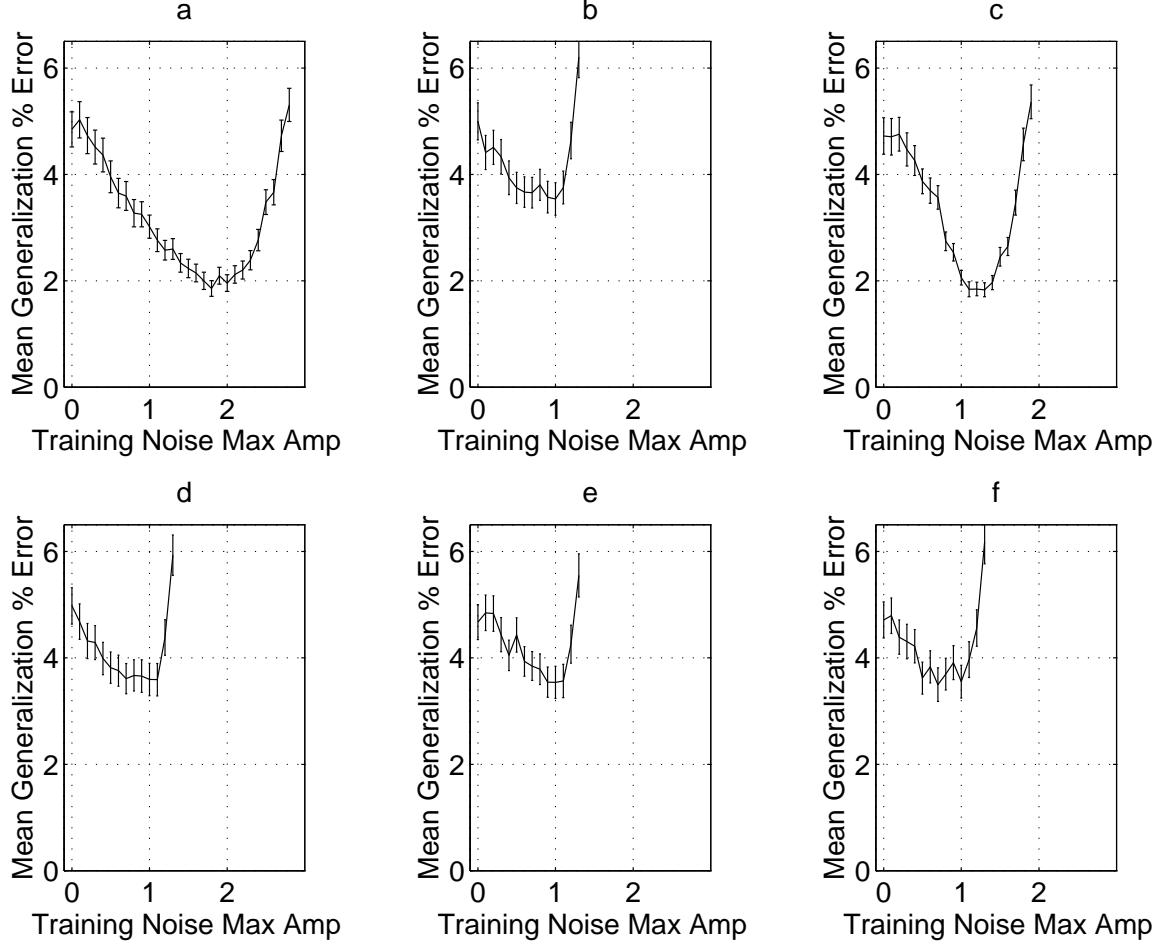


Figure 3: Mean generalization percent error as a function of noise amplitude a for recurrent networks with 4-state neurons. The error bars show 90% confidence intervals, obtained by running 500 simulations for each point. All 3 curves in the top row are for additive noise and all 3 in the bottom row for multiplicative noise. The generalization curves are for (a) non-cumulative additive per time step; (b) non-cumulative additive per string; (c) cumulative additive per time step; (d) non-cumulative multiplicative per time step; (e) non-cumulative multiplicative per string; (f) cumulative multiplicative per time step. Experiments with cumulative per string noise do not converge.

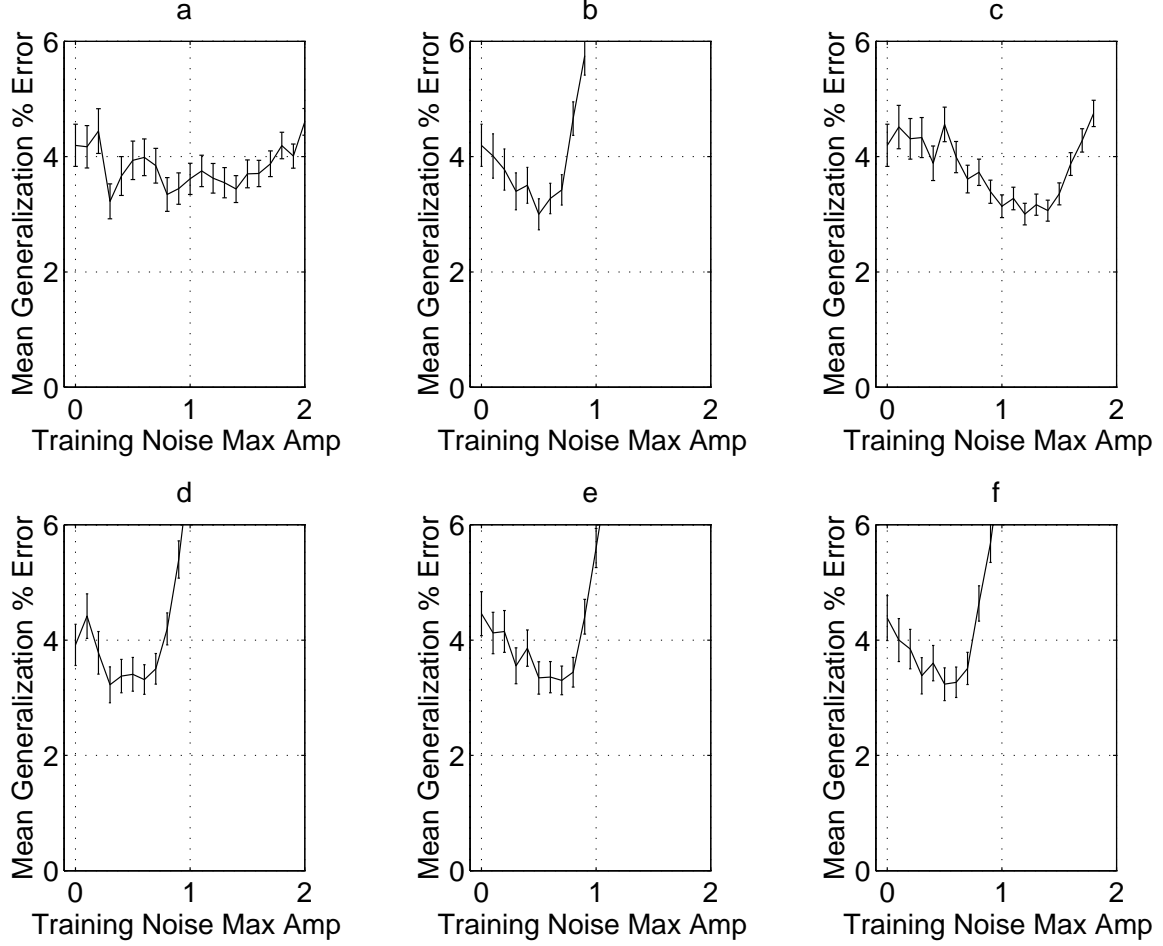


Figure 4: Mean generalization percent error as a function of error magnitude a for recurrent networks with 3-state neurons. The error bars show 90% confidence intervals, obtained by running 500 simulations for each data point. All 3 curves in the top row are for additive noise and all 3 in the bottom row for multiplicative noise. The generalization curves are for (a) non-cumulative additive per time step; (b) non-cumulative additive per string; (c) cumulative additive per time step; (d) non-cumulative multiplicative per time step; (e) non-cumulative multiplicative per string; (f) cumulative multiplicative per time step. Experiments with cumulative per string noise do not converge.

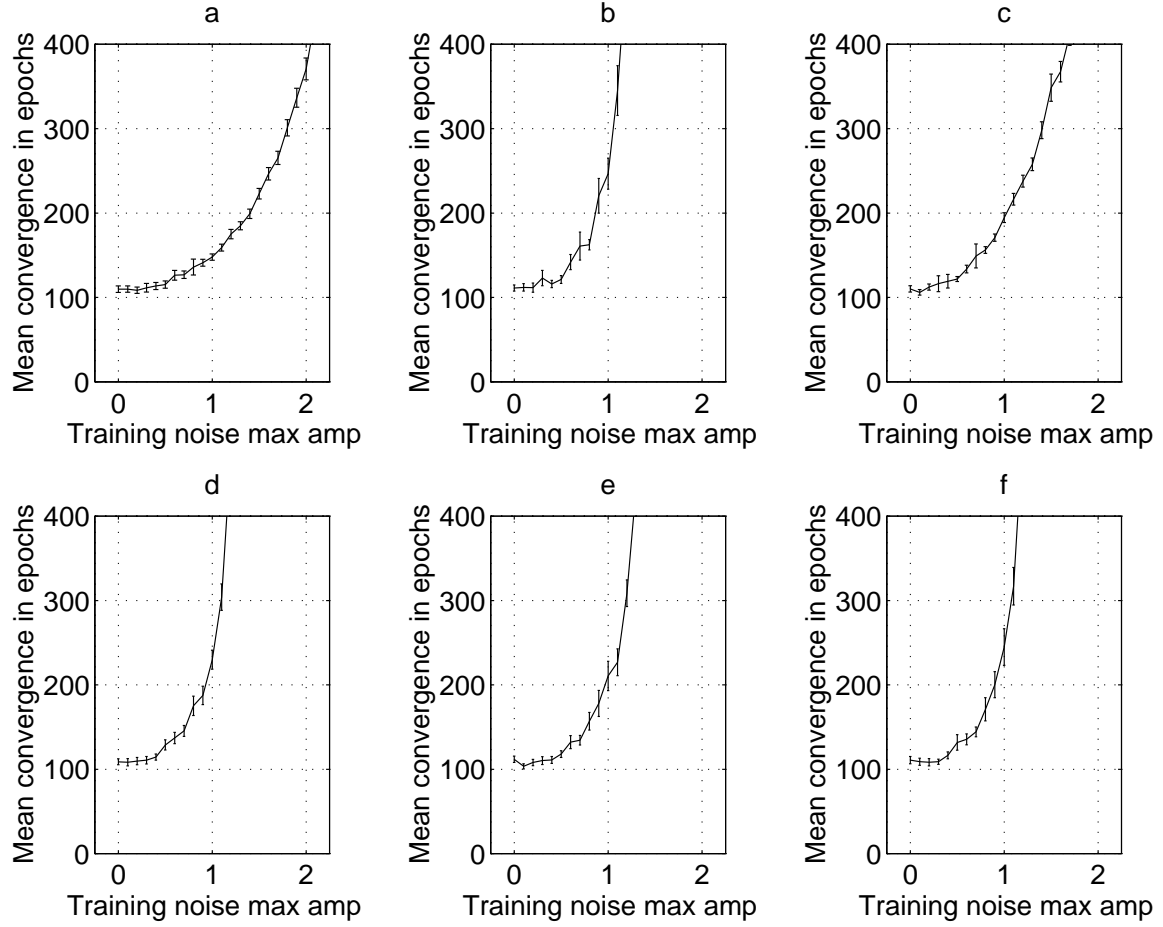


Figure 5: Mean convergence in epochs as a function of noise amplitude a for recurrent networks with 4-state neurons. The error bars show 90% confidence intervals, obtained by running 500 simulations for each data point. The convergence curves are for (a) non-cumulative additive per time step; (b) non-cumulative additive per string; (c) cumulative additive per time step; (d) non-cumulative multiplicative per time step; (e) non-cumulative multiplicative per string; (f) cumulative multiplicative per time step. Experiments with cumulative per string noise do not converge.

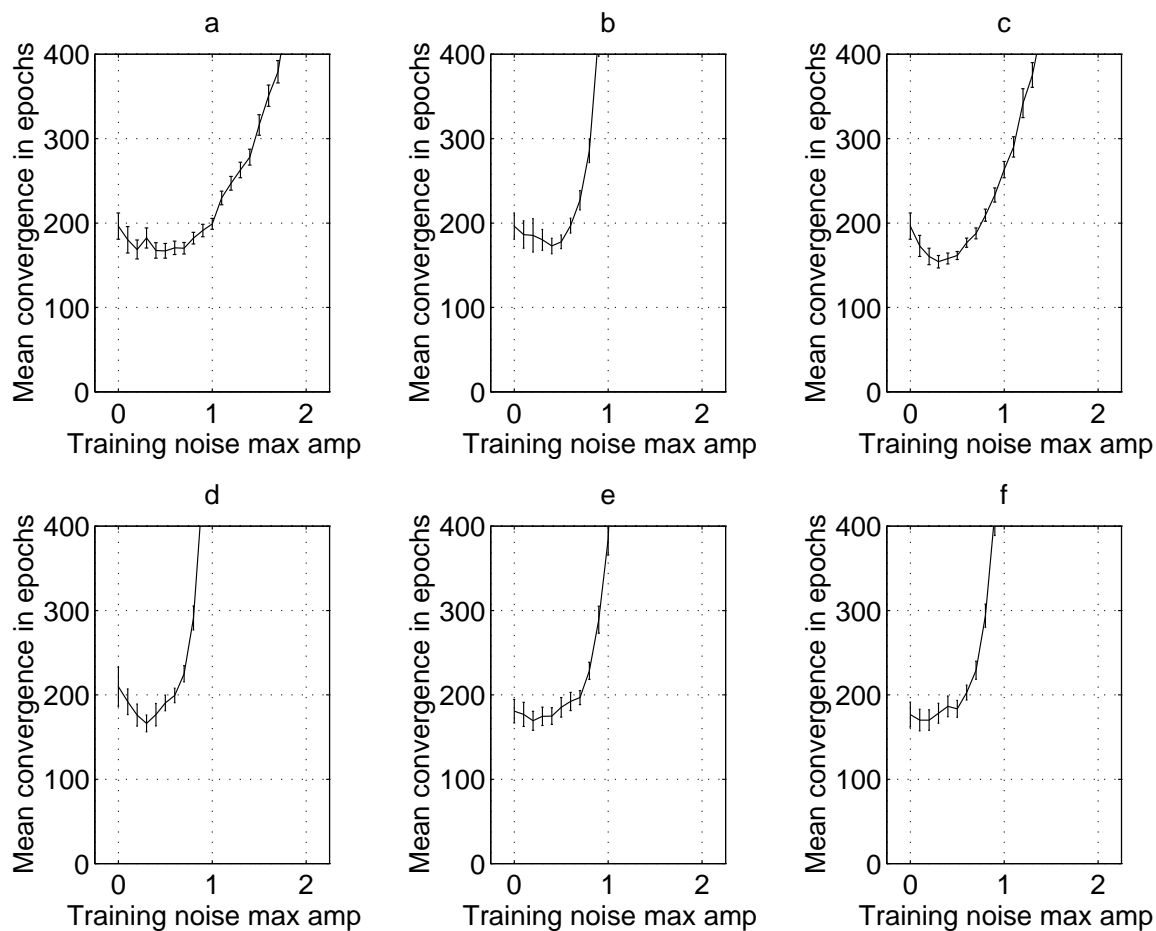


Figure 6: Mean convergence in epochs as a function of noise amplitude a for recurrent networks with 3-state neurons. The error bars show 90% confidence intervals, obtained by running 500 simulations for each data point. The convergence curves are for (a) non-cumulative additive per time step; (b) non-cumulative additive per string; (c) cumulative additive per time step; (d) non-cumulative multiplicative per time step; (e) non-cumulative multiplicative per string; (f) cumulative multiplicative per time step. Experiments with cumulative per string noise do not converge.

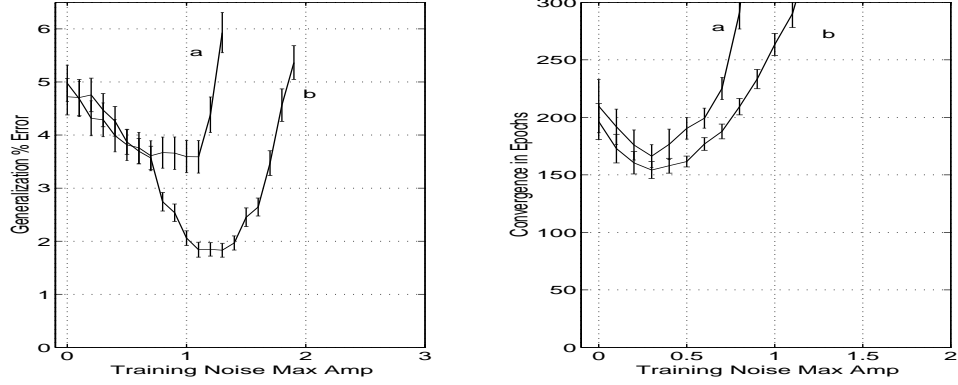


Figure 7: Best Convergence/Generalization for Additive and Multiplicative Noises. (a) multiplicative non-cumulative per time step; (b) additive cumulative per time step.

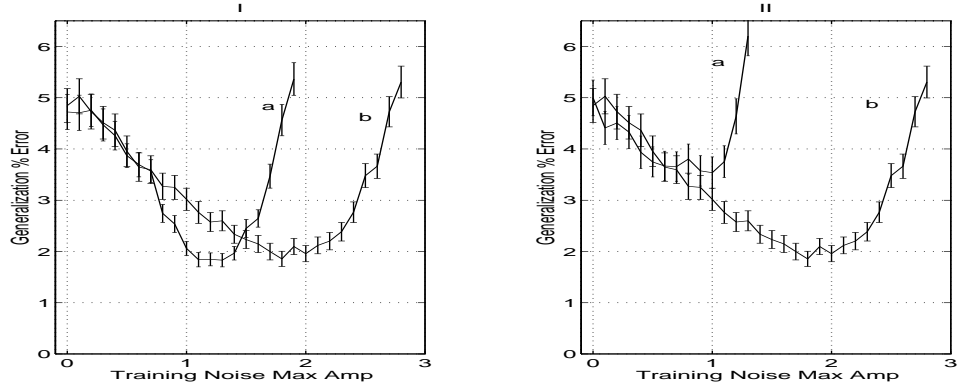


Figure 8: (I) Best Generalization for Cumulative and Non-Cumulative Noises: a) cumulative additive per time step; b) non-cumulative additive per time step. (II) Best Generalization for non-cumulative Per Time Step and Per String Noises: a) per string additive; b) per time step additive.

4.4 Effects of each Noise Parameter on Performance

This section makes use of the dual parity simulations to assess the effects of each noise parameter on convergence time and generalization performance. In order to make fair assessments, it is reasonable to compare the best mean performance due to each noise parameter. For example, to compare the effects of additive versus multiplicative noise on generalization, we should compare the mean generalization curves of networks injected with additive cumulative per time step noise with the mean generalization curves of networks injected with multiplicative non-cumulative per time step noise. These noises represent the additive and multiplicative noise combinations yielding the best generalization. By looking at the best noise combinations for each parameter, comparisons can be made which are independent of other noise issues. Using the above case as an example, if the best additive noise show better generalization than the best multiplicative noise, then this is a fair indication that additive noise yields better generalization than multiplicative noise. It cannot be argued that perhaps cumulative noise allows better generalization, since

if this was the case then we would be comparing additive cumulative per time step noise with multiplicative cumulative per time step noise – which would be fair if such was the case.

Simulated performance closely mirror our predictions. Improvements were observed for all noise models except for cumulative per string noises which failed to converge for all runs. Generalization improvement was more emphasized on networks with 4 states, while convergence enhancement was more noticeable on 3-state networks. The simulations show the following results:

- Additive noise is better tolerated than multiplicative noise, achieves better generalization, and yields slightly better convergence (Figure 7).
- Cumulative noise achieves optimal generalization and convergence at lower amplitudes than non-cumulative noise. Cumulative noise also has a narrower range of beneficial noise, which is defined as the range of noise amplitudes which yields better performance than that of a noiseless network (Figure 8a illustrates this for generalization).
- All non-cumulative per time step noises result in better convergence than non-cumulative per string noises on 3-state networks. Non-cumulative additive per time step noise yields better generalization and is better tolerated than non-cumulative additive per string noise (Figure 8b).

Overall, the best performance is obtained by applying cumulative and non-cumulative additive noise at each time step. These results closely match the predictions of section 3.1. The only exceptions are that all multiplicative noise models seem to yield equivalent performance. This discrepancy between prediction and simulation may be due to the detrimental effects of weight decay in multiplicative noise, which can conflict with the advantages of cumulative and per time step noise.

4.5 Generalization vs. Convergence

In most applications both generalization and convergence are important; improving one at the expense of the other becomes a difficult decision. In this section we plot the trade-off between generalization and convergence in the dual parity results, and show that in many cases it is possible to improve *both* performance measures simultaneously.

In figure 9 we plot generalization vs. convergence. Noise magnitude increases from left to right in each curve. For example, if convergence is not an issue than one would pick the curve that dips the lowest, and choose the corresponding noise amplitude for the noise method corresponding to that curve. These plots clearly illustrate that additive noise methods allow the best generalization and convergence for the 4-state neurons networks.

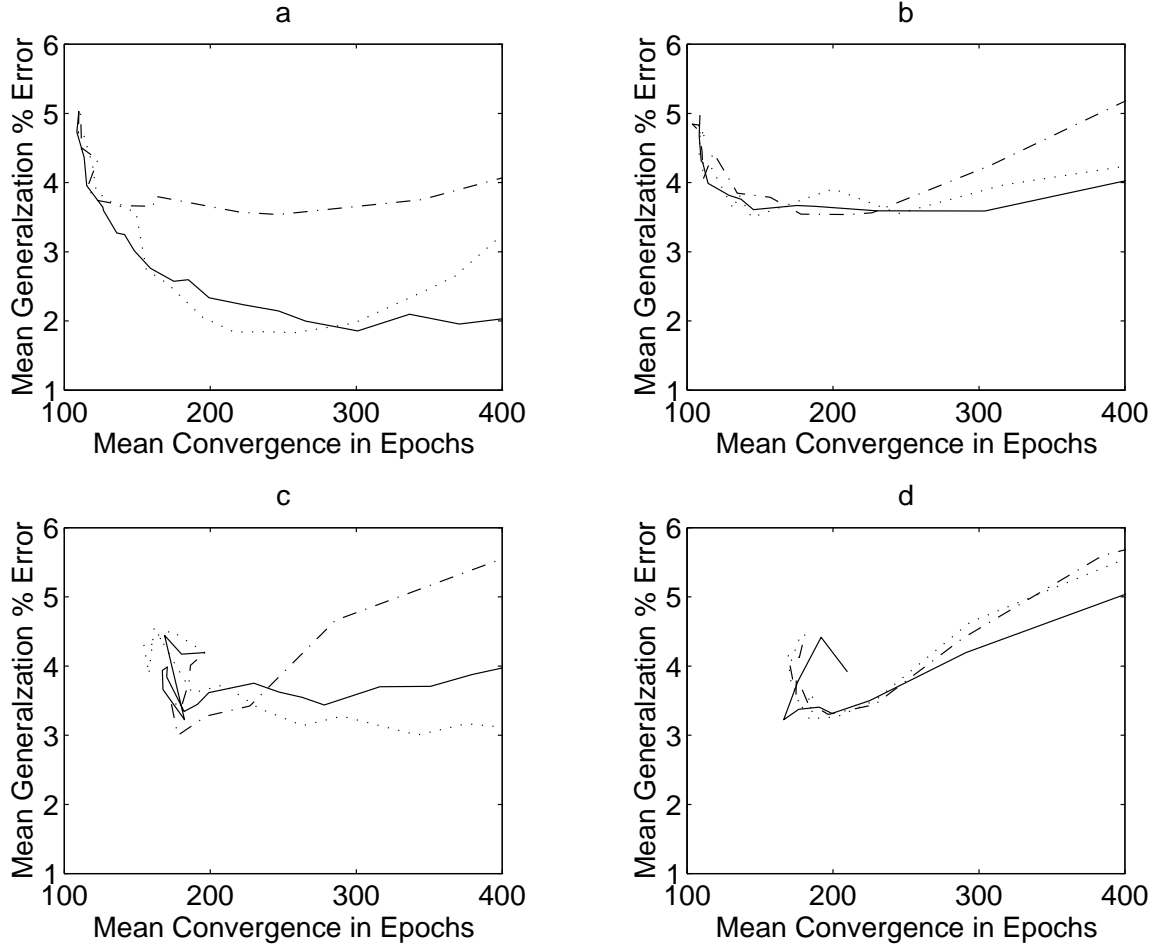


Figure 9: Generalization percent error rate vs. convergence in epochs. The curves are obtained by using the data in previous figures (recall 500 simulations for each data point). The lines represent the following: – non-cumulative per time step; - - non-cumulative per string; ... cumulative per time step. The plots represent the following: (a) additive noise on 4-state network; (b) multiplicative noise on 4-state network; (c) additive noise on 3-state network; (d) multiplicative noise on 3-state network. Noise magnitude increases from left to right in each curve.

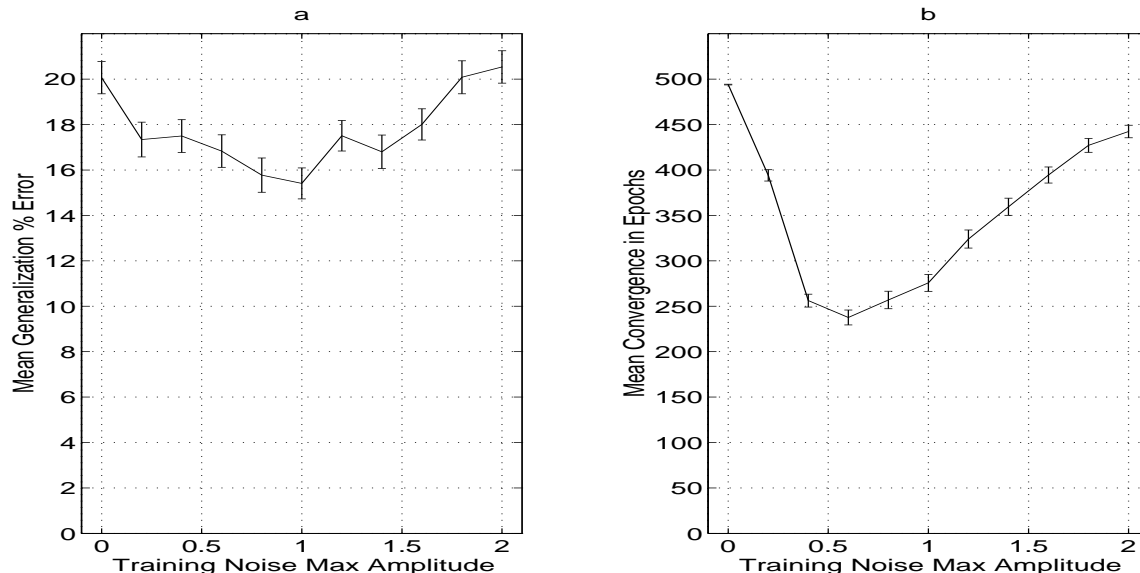


Figure 10: Generalization and convergence results on the RND6 grammar. Additive, non-cumulative, per time step noise is inserted into second order recurrent networks with 10 state neurons. (a) Mean generalization percent error as a function of noise amplitude. (b) Mean convergence in epochs as a function of noise amplitude. The error bars show 90% confidence intervals, obtained by running 500 simulations for each point.

More importantly, these plots illustrate the cases where both convergence and generalization are improved. In figures 9c and 9d the curves clearly curl down and to the left for lower noise amplitudes before rising to the right at higher noise amplitudes. These lower regions are optimal because they represent improved generalization and improved convergence, i.e., by picking a noise amplitude from this “hook” convergence and generalization do not trade-off.

4.6 Convergence and Generalization Results on the RND6 Grammar

Figures 10 and 11 plot convergence and generalization results on the RND6 grammar using additive, non-cumulative, per time step noise. These results are obtained from a problem slightly more difficult than the dual parity problem, in the sense that the RND6 grammar is slightly larger and no attempts have been made to optimize the network size for this grammar. A network size of 10 state neurons was chosen arbitrarily, although more neurons could have been added to improve generalization performance.

We simulated only the best predicted noise model for this set of experiments. The purpose was to observe the performance on a slightly more difficult problem. As shown in Figures 10 and 11, the results are consistent with those on the dual parity grammar presented above. Both generalization and convergence improved for a range of noise amplitudes.

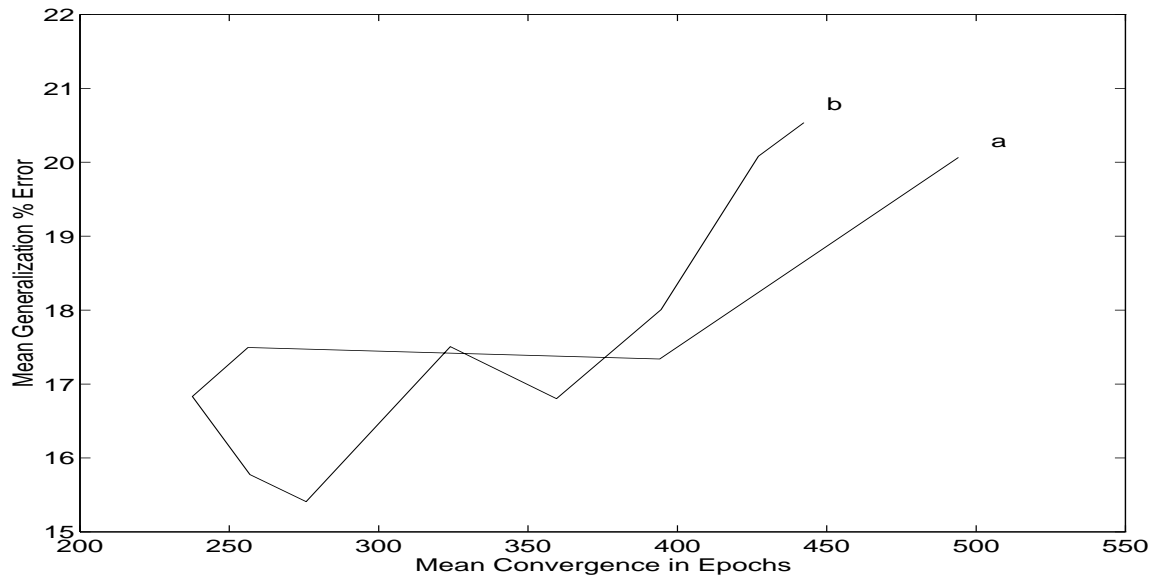


Figure 11: Generalization percent error vs. convergence in epochs for the RND6 grammar. Additive, non-cumulative, per time step noise is inserted into second order recurrent networks with 10 state neurons. Maximum noise magnitude m increases from point a ($m = 0$) to point b ($m = 2.0$).

5 CONCLUSIONS

We have presented several methods of injecting synaptic noise to recurrent neural networks. We give a thorough analysis of these methods and empirically test these methods on the learning dual parity and RND6 grammars from strings encoded as temporal sequences.

Results show that most of these methods offer several benefits on recurrent networks:

- Ability to improve generalization performance.
- Ability to improve convergence performance.
- Ability to improve generalization and convergence simultaneously.

These benefits should apply to feedforward and other recurrent networks. The ability to improve generalization and convergence simultaneously makes these results novel to other methods previously discussed for feedforward nets, which in general cast convergence as a cost for improved generalization performance.

These results also offer some insight into what types of synaptic noises are better tolerated when training under noisy environments. If only synaptic noise is applicable, and noise-free updates are allowed, then cumulative multiplicative noise is the most detrimental at high amplitudes. In addition, noise which is uncorrelated in time is less destructive, as illustrated by the fact that injecting a different noise at each time step (for one input string) allows higher detrimental noise thresholds than injecting the same noise at each

time step (effectively equivalent to the per string noise models). However, the appropriateness of any of these noise environments needs to be determined.

It would be interesting to apply these methods to other problems to study the effect of problem complexity on the range of beneficial noise. Further research with larger networks could determine the effect of network size on noise tolerance and effectiveness on generalization and convergence. Another important issue is the relationship between the range of optimal noise values to network size and problem complexity. Thus it could be possible to determine the usefulness of synaptic noise by allowing an optimal range of noise values to be predicted for each application.

6 ACKNOWLEDGEMENTS

We would like to thank and acknowledge the reviewers for valuable comments and suggestions, and acknowledge Li Chun An for help in the simulations and useful discussions.

APPENDIX

A Series Expansions of the Error Functions

The Taylor series expansion [8, 24] terms of the error functions of all noise models are derived in this section. In all derivations below the noise is assumed to be zero-mean white and uncorrelated in time.

A.1 Non-Cumulative Noise

The additive, per time step noise models are analyzed in detail below, followed by straightforward extensions to other non-cumulative noise cases.

A.1.1 Additive Per Time Step

Each weight $W_{t,ijk}$ is perturbed at each time step such that

$$W_{t,ijk} \longrightarrow W_{ijk}^* + \Delta_{t,ijk} \quad (11)$$

where $\{W_{ijk}^*\}$ is the noise-free weight set. Performing a Taylor series expansion on the output of state neuron S_O^T to second order around the noise-free weight set yields

$$\begin{aligned} S_O^T \longrightarrow & S_O^T + \sum_{t=0}^{T-1} \sum_{ijk} \Delta_{t,ijk} \left(\frac{\partial S_O^T}{\partial W_{t,ijk}} \right) + \\ & \frac{1}{2} \sum_{t,u=0}^{T-1} \sum_{ijk,lmn} \Delta_{t,ijk} \Delta_{u,lmn} \left(\frac{\partial^2 S_O^T}{\partial W_{t,ijk} \partial W_{u,lmn}} \right) + O(>3) \end{aligned} \quad (12)$$

The error E_p for pattern p on a noiseless network is defined as

$$E_p = \frac{1}{2} \epsilon_p^2 \quad (13)$$

where ϵ_p is

$$\epsilon_p = S_O^T - d_p. \quad (14)$$

However, the injection of noise represented by equation (12) augments E_p as follows

$$E_p = \frac{1}{2} [S_O^T + \sum_{t=0}^{T-1} \sum_{ijk} \Delta_{t,ijk} \left(\frac{\partial S_O^T}{\partial W_{t,ijk}} \right) +$$

$$\frac{1}{2} \sum_{t,u=0}^{T-1} \sum_{ijk,lmn} \Delta_{t,ijk} \Delta_{u,lmn} \left(\frac{\partial^2 S_O^T}{\partial W_{t,ijk} \partial W_{u,lmn}} \right) + O(>3) - d_p]^2. \quad (15)$$

Dropping Δ^3 and higher order terms,

$$\begin{aligned} E_p = & \frac{1}{2} \left[\epsilon_p^2 + \sum_{t,u=0}^{T-1} \sum_{ijk,lmn} \Delta_{t,ijk} \Delta_{u,lmn} \left(\frac{\partial S_O^T}{\partial W_{t,ijk}} \right) \left(\frac{\partial S_O^T}{\partial W_{u,lmn}} \right) + \right. \\ & 2\epsilon_p \sum_{t=0}^{T-1} \sum_{ijk} \Delta_{t,ijk} \left(\frac{\partial S_O^T}{\partial W_{t,ijk}} \right) + \\ & \left. \epsilon_p \sum_{t,u=0}^{T-1} \sum_{ijk,lmn} \Delta_{t,ijk} \Delta_{u,lmn} \left(\frac{\partial^2 S_O^T}{\partial W_{t,ijk} \partial W_{u,lmn}} \right) \right]. \quad (16) \end{aligned}$$

Since the noise is zero-mean white and uncorrelated, taking the time average over the training phase for the per time step case yields

$$\langle \Delta_{t,ijk} \rangle = 0 \quad (17)$$

$$\langle \Delta_{t,ijk} \Delta_{u,lmn} \rangle = \sigma^2 \delta_{tijk,ulmn} \quad (18)$$

where $\delta_{tijk,ulmn}$ is the Kronecker delta. Using equations (17)(18) simplifies (16) to

$$\langle E_p^{noisy} \rangle = \langle E_p \rangle + \frac{1}{2} \sum_{t=0}^{T-1} \sum_{ijk} \sigma^2 \left[\left(\frac{\partial S_O^T}{\partial W_{t,ijk}} \right)^2 + \epsilon_p \left(\frac{\partial^2 S_O^T}{\partial (W_{t,ijk})^2} \right) \right]. \quad (19)$$

It is important to realize that although each “layer” has the same set of weights,

$$W_{0,ijk} = W_{1,ijk} = \dots = W_{T-1,ijk} \quad (20)$$

nevertheless their respective partial derivatives are not equal because they are time dependent:

$$\frac{\partial S_O^T}{\partial W_{0,ijk}} \neq \frac{\partial S_O^T}{\partial W_{1,ijk}} \neq \dots \neq \frac{\partial S_O^T}{\partial W_{T-1,ijk}}. \quad (21)$$

Therefore all partial derivatives $\sum_{t=0}^{T-1} \sum_{ijk} [(\frac{\partial S_O^T}{\partial W_{t,ijk}})^2 + (\frac{\partial^2 S_O^T}{\partial (W_{t,ijk})^2})]$ need to be accounted for in equation (19).

A.1.2 Additive Per String

To extend the above derivation to the per string case, note that $\Delta_{t,ijk} = \Delta_{ijk}$ for all time steps within each string in eq. (11). Therefore,

$$\begin{aligned} \langle \Delta_{t,ijk} \rangle &= 0 \\ \langle \Delta_{t,ijk} \Delta_{u,lmn} \rangle &= \sigma^2 \delta_{ijk,lmn} \end{aligned} \quad (22)$$

and eq. (19) becomes

$$\langle E_p^{noisy} \rangle = \langle E_p \rangle + \frac{1}{2} \sum_{t,u=0}^{T-1} \sum_{ijk} (\sigma^2) \left[\left(\frac{\partial S_O^T}{\partial W_{t,ijk}} \right) \left(\frac{\partial S_O^T}{\partial W_{u,ijk}} \right) + \epsilon_p \left(\frac{\partial^2 S_O^T}{\partial W_{t,ijk} \partial W_{u,ijk}} \right) \right]. \quad (23)$$

A.1.3 Multiplicative Per Time Step and Per String

The above per string and per time step derivations can be easily extended to the multiplicative noise models by substituting

$$W_{t,ijk} \longrightarrow W_{ijk}^* + W_{ijk}^* \Delta_{t,ijk}. \quad (24)$$

in place of eq. (24).

A.2 Cumulative Noise

The cumulative noise models are analyzed below. They are based on and described in the same format as the above section on non-cumulative noise.

A.2.1 Additive Per Time Step

In the cumulative additive per time step case, each noise injection step is such that

$$W_{t,ijk} = W_{ijk}^* + \sum_{\tau=0}^t \Delta_{\tau,ijk}. \quad (25)$$

The time averages

$$\langle \sum_{\tau=0}^t \Delta_{\tau,ijk} \rangle = 0 \quad (26)$$

$$\langle \sum_{\tau=0}^t \Delta_{\tau,ijk} \sum_{\gamma=0}^u \Delta_{\gamma,lmn} \rangle = v \sigma^2 \delta_{ijk,lmn} \quad (27)$$

where $v = \min(t + 1, u + 1)$ can be used to simplify the error expression.

A.2.2 Additive Per String

The noise injection step for the additive per string case can be modeled by

$$W_{t,ijk} = W_{ijk}^* + (t + 1)\Delta_{ijk} \quad (28)$$

The time averages

$$\langle (t + 1)\Delta_{ijk} \rangle = 0 \quad (29)$$

$$\langle (t + 1)\Delta_{ijk}(u + 1)\Delta_{lmn} \rangle = (t + 1)(u + 1)\sigma^2\delta_{ijk,lmn} \quad (30)$$

can be used to simplify the error expression.

A.2.3 Multiplicative Per Time Step

For the multiplicative, per time step case, each noise injection step is such that

$$W_{t,ijk} = W_{ijk}^* \prod_{\tau=0}^t (1 + \Delta_{\tau,ijk}) \quad (31)$$

To put eq. (31) in a similar format for the error expansion, note that it can be represented by

$$W_{t,ijk} = W_{ijk}^* + W_{ijk}^* \Phi_{t,ijk} \quad (32)$$

where

$$\Phi_{t,ijk} = \prod_{\tau=0}^t (1 + \Delta_{\tau,ijk}) - 1. \quad (33)$$

The time averages

$$\langle \Phi_{t,ijk} \rangle = 0 \quad (34)$$

$$\langle \Phi_{t,ijk}\Phi_{u,lmn} \rangle = v\sigma^2\delta_{ijk,lmn} \quad (35)$$

where $v = \min(t + 1, u + 1)$ can be used to simplify the error expression.

A.2.4 Multiplicative Per String

In the per string case, the noise injection step $W_{t,ijk} = W_{ijk}^*(1 + \Delta_{ijk})^{t+1}$ can be represented by

$$W_{t,ijk} = W_{ijk}^* + W_{ijk}^* \Psi_{t,ijk} \quad (36)$$

where

$$\Psi_{t,ijk} = (1 + \Delta_{ijk})^{t+1} - 1 \quad (37)$$

Using binomial expansion and dropping third or higher terms we get the following approximations:

$$\begin{aligned} \langle \Psi_{t,ijk} \rangle &= \left[1 + (t+1) \langle \Delta_{ijk} \rangle + \frac{1}{2} t(t+1) \langle \Delta_{ijk}^2 \rangle + O(>3) \right] - 1 \\ &= \frac{1}{2} t(t+1) \sigma^2 \end{aligned} \quad (38)$$

$$\begin{aligned} \langle \Psi_{t,ijk} \Psi_{u,lmn} \rangle &= \langle (1 + \Delta_{ijk})^{t+1} (1 + \Delta_{lmn})^{u+1} - (1 + \Delta_{ijk})^{t+1} - (1 + \Delta_{lmn})^{u+1} + 1 \rangle \\ &= \langle (1 + \Delta_{ijk})^{t+1} (1 + \Delta_{lmn})^{u+1} \rangle - \frac{1}{2} t(t+1) \sigma^2 - \frac{1}{2} u(u+1) \sigma^2 - 1 \\ &= \left[1 + \frac{t(t+1) + u(u+1)}{2} \sigma^2 + (t+1)(u+1) \sigma^2 \delta_{ijk,lmn} \right] - \frac{1}{2} t(t+1) \sigma^2 - \frac{1}{2} u(u+1) \sigma^2 - 1 \\ &= (t+1)(u+1) \sigma^2 \delta_{ijk,lmn} \end{aligned} \quad (39)$$

For all cases the weight update ΔW_{ijk} descends the noise-enhanced error function $\langle E_p^{noisy} \rangle$.

B Cumulative vs. Non-cumulative Noise in Feedforward Networks

In feedforward networks, cumulative and non-cumulative noise models are identical if using per sample training, but can have different variance if batch training is used (see analysis below). The methods described in this paper update the set of noise-free weights and accumulate noise between weight updates, as opposed to the continuous accumulation of noise across weight updates. Per sample updating and batch updating are considered below.

B.1 Per Sample Updating

Since the length of each sample (in time) is one, noise cannot accumulate in the weights because we do not keep track of the previous noise values after each weight update. Thus, cumulative and non-cumulative noise

models are identical when updating the weights after the presentation of each training sample.

B.2 Batch Training

In this case, the errors for several samples are computed before the weights are updated. For example, in an additive noise model the weights are augmented as follows between updates, for each training sample t_i and weight W_j :

$$\begin{aligned} t_0 & : W_j = W_j + n_0, & n_0 &= \Delta_0 \\ t_1 & : W_j = W_j + n_1, & n_1 &= \Delta_0 + \Delta_1 \\ t_2 & : W_j = W_j + n_2, & n_2 &= \Delta_0 + \Delta_1 + \Delta_2 \\ & \vdots \\ t_M & : W_j = W_j + n_M, & n_M &= \Delta_0 + \Delta_1 + \Delta_2 + \cdots + \Delta_M \end{aligned}$$

where M is the number of training samples in the batch, Δ_i are uncorrelated noise values and n_i are the effective non-cumulative noise values. Although the terms n_i are correlated, the correlations do not contribute to the augmented error function in the feedforward case. (They do contribute in recurrent architectures, as shown in Appendix A).

Thus, inserting cumulative noise Δ_i between updates is effectively equivalent to inserting noncumulative noise values n_i between updates. The noise values n_i are also zero-mean, however, the variance increases linearly with the number of training samples in the batch. This is easily shown by using the following equality:

$$Var(aX + bY) = a^2 Var(X) + b^2 Var(Y) + 2ab Cov(X, Y) \quad (40)$$

where a and b are constants, X and Y are random variables, $Var()$ is the variance, and $Cov()$ is the covariance. Applying this equality here, X and Y are independent noise values, $a = 1$, $b = 1$, $Cov(X, Y) = 0$, and the following is obtained:

$$Var(n_M) = Var(\Delta_0) + Var(\Delta_1) + \cdots + Var(\Delta_M). \quad (41)$$

Thus, the noise values for training samples presented later within the batch will observe greater variance.

References

- [1] D.H. Ackley, G.E. Hinton, and T.J. Sejnowski. A learning algorithm for boltzmann machines. *Cognitive Science*, 9:147–169, 1985.
- [2] E.B. Baum and D. Haussler. What size net gives valid generalization? *Neural Computation*, 1:151–160, 1989.
- [3] C. M. Bishop. Training with noise is equivalent to Tikhonov regularization. *Neural Computation*, 7(1):108–116, 1995.
- [4] H. J. Breermann and R. W. Anderson. An alternative to back-propagation: A simple rule of synaptic modification for neural net training and memory. Technical Report PAM-483, U. C. Berkeley Center for Pure and Applied Mathematics, 1989.
- [5] H. J. Breermann and R. W. Anderson. How the brain adjusts synapses - maybe. In R. S. Boyer, editor, *Automated Reasoning: Essays in Honor of Woody Bledsoe*, chapter 6, pages 119–147. Kluwer Academic Publishers, Boston, MA, 1991.
- [6] Robert M. Burton, Jr. and George J. Mpitsos. Event-dependent control of noise enhances learning in neural networks. *Neural Networks*, 5:627–637, 1992.
- [7] Gert Cauwenberghs. A fast stochastic error-descent algorithm for supervised learning and optimization. In S.J. Hanson, J.D. Cowan, and C.L. Giles, editors, *Advances in Neural Information Processing Systems 5*, pages 244–251, San Mateo, CA, 1993. Morgan Kaufmann Publishers.
- [8] Yann Le Cun, John S. Denker, and Sara A. Solla. Optimal brain damage. In D. Touretzky, editor, *Advances in Neural Information Processing Systems 2*, pages 598–6053. Morgan Kaufmann Publishers, 1990.
- [9] A. Dembo and T. Kailath. Model-free distributed learning. *IEEE Transactions on Neural Networks*, 1:58–70, 1990.
- [10] Barry Flower and Marwan Jabri. Summed weight neuron perturbation: An $O(n)$ improvement over weight perturbation. In S.J. Hanson, J.D. Cowan, and C.L. Giles, editors, *Advances in Neural Information Processing Systems 5*, pages 212–219, San Mateo, CA, 1993. Morgan Kaufmann Publishers.
- [11] C.L. Giles, C.B. Miller, D. Chen, H.H. Chen, G.Z. Sun, and Y.C. Lee. Learning and extracting finite state automata with second-order recurrent neural networks. *Neural Computation*, 4(3):393–405, 1992.
- [12] C.L. Giles and C.W. Omlin. Pruning recurrent neural networks for improved generalization performance. *IEEE Transactions on Neural Networks*, 5(5):848–851, 1994.
- [13] Stephen José Hanson. A stochastic version of the delta rule. *Physica D.*, 42:265–272, 1990.
- [14] G.E. Hinton and T.J. Sejnowski. Learning and relearning in Boltzmann machines. In *Parallel Distributed Processing*, chapter 7, pages 282–317. MIT Press, Cambridge, MA, 1986.
- [15] Geoffrey E. Hinton, Terrence J. Sejnowski, and David H. Ackley. Boltzmann machines: Constraint satisfaction networks that learn. Technical Report CMU-CS-84-119, Carnegie-Mellon University CS Dept., May 1984.
- [16] Lasse Holmström and Petri Koistinen. Using additive noise in back-propagation training. *IEEE Transactions on Neural Networks*, 3(1):24–38, 1992.
- [17] Marwan Jabri and Barry Flower. Weight perturbation: An optimal architecture and learning technique for analog vlsi feedforward and recurrent multilayer networks. *IEEE Transactions on Neural Networks*, pages 154–157, 1992.

- [18] Stephen Judd and Paul W. Munro. Nets with unreliable hidden nodes learn error-correcting codes. In S.J. Hanson, J.D. Cowan, and C.L. Giles, editors, *Advances in Neural Information Processing Systems 5*, pages 89–96, San Mateo, CA, 1993. Morgan Kaufmann Publishers.
- [19] Danny Kilis and Eugene Ackerman. A stochastic neuron model for pattern recognition. In *Proceedings of the International Joint Conference on Neural Networks*, volume 2, Washington, D.C., January 1990.
- [20] S. Kirkpatrick, Jr. C.D. Galett, and M.P. Vecchi. Optimization by simulated annealing. *Science*, 220:671–680, May 1983.
- [21] Anders Krogh and John A. Hertz. A simple weight decay can improve generalization. In J.E. Moody, S.J. Hanson, and R.P. Lippmann, editors, *Advances in Neural Information Processing Systems 4*, pages 450–957, San Mateo, CA, 1992. Morgan Kaufmann Publishers.
- [22] C.B. Miller and C.L. Giles. Experimental comparison of the effect of order in recurrent neural networks. *International Journal of Pattern Recognition and Artificial Intelligence*, 7(4):849–872, 1993. Special Issue on Neural Networks and Pattern Recognition, editors: I. Guyon , P.S.P. Wang.
- [23] Jay. I. Minnix. Fault tolerance of the backpropagation neural network trained on noisy inputs. In *Proc. of IJCNN*, volume I, pages I–847–852, 1992.
- [24] Alan F. Murray and Peter J. Edwards. Enhanced MLP performance and fault tolerance resulting from synaptic weight noise during training. *IEEE Trans. on Neural Networks*, 5(5):792–802, 1994.
- [25] C.W. Omlin and C.L. Giles. Fault-tolerant implementation of finite-state automata in recurrent neural networks. Technical Report 95-3, Computer Science Department, Rensselaer Polytechnic Institute, Troy, NY 12180, 1995.
- [26] Russel Reed. Pruning algorithms – a survey. *IEEE Trans. on Neural Networks*, 4(5):740–747, 1993.
- [27] R.J. Williams and D. Zipser. Gradient-based learning algorithms for recurrent networks and their computational complexity. In Y. Chauvin and D. E. Rumelhart, editors, *Back-propagation: Theory, Architectures and Applications*, chapter 13, pages 433–486. Lawrence Erlbaum Publishers, Hillsdale, N.J., 1995.
- [28] Carlo H. Séquin and Reed D. Clay. Fault tolerance in artificial neural networks. In *Proc. of IJCNN*, volume I, pages I–703–708, 1990.
- [29] R.L. Watrous and G.M. Kuhn. Induction of finite-state languages using second-order recurrent networks. *Neural Computation*, 4(3):406, 1992.
- [30] Z. Zeng, R.M. Goodman, and P. Smyth. Learning finite state machines with self-clustering recurrent networks. *Neural Computation*, 5(6):976–990, 1993.