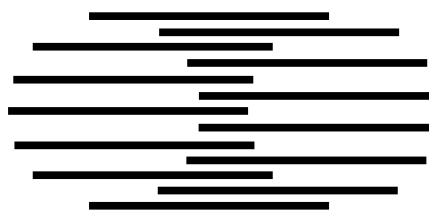


IDIAP  
Technical report



## High Order and Multilayer Perceptron Initialization

G. Thimm and E. Fiesler

*November, 1994*

This paper is also submitted to  
*IEEE Transactions on Neural  
Networks*

## Abstract

Proper initialization is one of the most important prerequisites for fast convergence of feed-forward neural networks like high order and multilayer perceptrons. This publication aims at determining the optimal value of the initial weight variance (or range), which is the principal parameter of random weight initialization methods for both types of neural networks.

An overview of random weight initialization methods for multilayer perceptrons is presented. These methods are extensively tested using eight real-world benchmark data sets and a broad range of initial weight variances by means of more than 30,000 simulations, in the aim to find the best weight initialization method for multilayer perceptrons.

For high order networks, a large number of experiments (more than 200,000 simulations) was performed, using three weight distributions, three activation functions, several network orders, and the same eight data sets. The results of these experiments are compared to weight initialization techniques for multilayer perceptrons, which leads to the proposal of a suitable weight initialization method for high order perceptrons.

The conclusions on the weight initialization methods for both types of networks are justified by sufficiently small confidence intervals of the mean convergence times.

**Keywords:** *neural network initialization, random weight initialization, initial weight, weight initialization, interconnection strength, comparison of weight initialization methods, high(er) order neural network, high(er) order perceptron, sigma-pi connection, initial weight distribution, activation function, learning rate, multilayer perceptron, neural computation, neural network, neurocomputing, optimization, connectionism, real-world benchmark*

## 1 Introduction

The learning speed of multilayer and high order perceptrons<sup>1</sup> depends mainly on the initial values of its weights and biases, its learning rate, its network topology, and on learning rule improvements like the momentum term. The optimal values for these parameters are usually unknown *a priori* because they depend mainly on the training data set used. In practice it is not feasible to perform a global search for obtaining the optimal values of these parameters, as the convergence behavior of the network might change significantly for small changes in the initial weights, as was demonstrated by J. F. Kolen and J. B. Pollack [Kolen-90]. An extensive search for the optimum values requires therefore much more overhead than performing a relatively small number of simulations using non-optimal values. Furthermore, current mathematical techniques are insufficient for a complete theoretical study of the learning behavior of these neural networks. Nevertheless, it is important to have a good approximation of the optimal initial value of these parameters; or with the words of J. F. Kolen and J. B. Pollack: to start the learning process in the “eye of the storm,” to reduce the required training time.

Several weight initialization methods for multilayer perceptrons have been suggested. The simplest method among them is random weight initialization, which is often preferred for its simplicity and its ability to produce multiple solutions, as the weights may, due to their initial randomness, converge to various attractors [Kolen-90]. Other methods involve extensive statistical and/or geometrical analysis of the data and are therefore very time consuming. The most rigorous among those is the pseudo-inverse method for perceptrons, which, besides being limited to linear separable data, has several other drawbacks (see [Hertz-91]). Some other weight initialization methods are based on special properties of a network that can not be applied to high order or multilayer perceptrons, as for example the weight initialization technique for radial basis function networks by J. C. Platt [Platt-91].

D. E. Rumelhart et al. observed that if all weights in a neural network are initialized with zero, they have the tendency to assume identical values during training. They therefore proposed random weight initialization to avoid this undesired situation by *breaking the symmetry* [Rumelhart-86]. However, the efficiency of this method depends much on the initial weight distribution. Several researchers therefore proposed random weight initialization methods. An overview of these methods is presented in section 2, and their performance is evaluated in section 4.2.1.

---

<sup>1</sup> For a definition of high order neural networks and associated references see [Fiesler-94] and [Thimm-94.1].

In order to obtain a thorough insight in the initialization characteristics of high order networks, which have not been studied before, numerous experiments were performed, varying the following parameters:

- the shape of the initial weight distribution: uniform, normal or Gaussian<sup>2</sup>, and a novel distribution which is uniform over the intervals  $[-2a, -a]$  and  $[a, 2a]$ , and zero everywhere else,
- the variance of the initial weight distribution,
- the order and topology of the network, and
- the activation function.

The results of these experiments is a simple weight initialization method using an application independent variance. In section 5, this method is compared to methods developed for multilayer perceptrons, in order to profit from experiences of other researchers and to determine a best method for higher order perceptrons.

## 2 Weight Initialization Techniques for Multilayer Perceptrons

S. E. Fahlman performed an early experimental study on the random weight initialization scheme for multilayer perceptrons. Based on this study, he proposed to use a uniform distribution with a range of  $[-1.0, 1.0]$ , but found that the best weight range for the data sets in his study varied from  $[-4.0, 4.0]$  to  $[-0.5, 0.5]$  [Fahlman-88].

Other researchers tried to determine the optimal weight range using network parameters:

L. Bottou uses an interval  $[-a/\sqrt{d_{in}}, a/\sqrt{d_{in}}]$ , where  $a$  is chosen in such way that the weight variance corresponds to the points of the maximal curvature of the activation function (which is approximately 2.38 for a standard sigmoid), and  $d_{in}$  is the fan-in (or in-degree) of a neuron, without justifying this interval further in a theoretical manner. L. Bottou trains the neural network only on speech data and does not compare this method with others [Bottou-88].

J. W. Boers and H. Kuiper initialize weights using a uniform distribution over the interval  $[-3/\sqrt{d_{in}}, 3/\sqrt{d_{in}}]$ , without any mathematical justification. They state that this interval performed the best on their speech data [Boers-92].

F. J. Śmieja uses uniformly distributed weights which are normalized to the magnitude  $2/\sqrt{d_{in}}$  for each node. The thresholds of the hidden units are then initialized to a random value in the interval  $[-\sqrt{d_{in}}/2, \sqrt{d_{in}}/2]$  and the thresholds of the output nodes are set to zero. He obtained these values from reasoning about hyperplane spin dynamics, and did not validate his method by experiments [Śmieja-91].

L. F. A. Wessels and E. Barnard describe two weight initialization methods. The first method sets the initial weight range to a value which assumes that the output of the network and the output patterns have the same variance. The second method puts equally distributed decision boundaries in the input space (without considering input or output patterns), which produces initial weights for the first interlayer weight matrix. The weights of the second interlayer weight matrix are set to 1.0. They compared both methods on generalization for three data sets. They found that the second method performed better in terms of generalization. However, they did not compare convergence speeds [Wessels-92].

An approach similar to the first method of L. F. A. Wessels and E. Barnard was introduced by G. P. Drago and S. Ridella [Drago-92]. They aim at avoiding flat regions in the error surface by restricting the number of neurons with absolute activations greater than 0.9. They developed a simple formula to estimate the best weight initialization scheme for multilayer perceptrons and showed for three data sets that this scheme uses satisfactory good initial weight ranges. The weights are uniformly distributed over the interval  $[-a, a]$ , with  $a = 1.3/\sqrt{1 + E[x^2]}$  for the input layer and  $a = 1.3/\sqrt{1 + 0.3d_{in}}$  for the output layer (assuming that all input values  $x$  have the same expected value  $E$ ).

Y. Lee, S.-H. Oh, and M. W. Kim showed theoretically that the probability of prematurely saturated neurons (small weight changes cause only negligible changes of the neuron output) in multilayer perceptrons increases with the maximal value of weights. They conclude that a smaller initial weight

---

<sup>2</sup>Neural network weights are often assumed to be normally distributed [Bellido-93].

range increases the learning speed of multilayer perceptrons. Simulations performed using two data sets confirm their reasoning, but they disregard that learning speed also decreases for weight ranges that are too small. Y. Lee et al. do not suggest an optimal weight range [Lee-93].

P. Haffner, A. Waibel, H. Sawai, and K. Shikano use a normal initial weight distribution. Unfortunately they do not compare their approach to others, give details, or justify it mathematically [Haffner-88].

R. L. Watrous and G. M. Kuhn compared a Gaussian distribution to a uniform distribution and found differences on the conditioning of the Jacobian matrix of a neural network, but found no relation to the convergence speed [Watrous-93].

D. Nguyen and B. Widrow use a multilayer perceptron with piecewise linear activation functions as an approximation of a network with logistic activation functions. Based on this simplification, they calculated an optimal length of  $d_i\sqrt{N_2}$  for the randomly initialized weight vectors and an optimal bias range of  $[-d_i\sqrt{N_2}, d_i\sqrt{N_2}]$  for neurons in the hidden layer, where  $N_2$  is the number of hidden nodes. The weights of the neurons in the output layer are randomly initialized in the interval  $[-0.5, 0.5]$ , without any justification given [Nguyen-90].

Y. K. Kim and J. B. Ra calculated a lower bound for the initial length of the weight vector of a neuron to be  $\sqrt{\eta/d_{in}}$ , where  $\eta$  is the learning rate [Kim-91].

Besides these random weight initialization methods, some non-random methods are described here for completeness.

A mixture between a random weight initialization scheme and the pseudo inverse method was developed by C.-L. Chen and R. S. Nutter for perceptrons with one hidden layer. First, the weights in the first interlayer weight matrix of the network are initialized with random values. Then, the weights in the second interlayer weight matrix are calculated using the pseudo inverse method applied to the activation values of the hidden layer. C.-L. Chen et al. refined this technique further by alternating the adjustment of the first interlayer weight matrix in a backpropagation-like process with the mentioned method of calculating the second interlayer weight matrix. These adjustments are repeated until a convergence criterion is reached, after which the backpropagation training begins. The authors report faster training in number of backpropagation cycles [Chen-91], but they disregard the computational complexity of the matrix inversions.

T. Denoeux and R. Lengellé initialize a one hidden layer perceptron with prototypes. This method requires a transformation of the input patterns to vectors of unit length and increased size. Additionally, prototypes have to be found by a cluster analysis. The authors reported improvements in training time, robustness versus local minima, and better generalization [Denoeux-93].

### 3 High Order Perceptrons

High order perceptrons are high order neural networks [Lee-86], having only unidirectional interlayer connections. They are also a generalization of Sigma-Pi networks, which are multilayer perceptrons having high order connections [Rumelhart-86], and functional link networks [Pao-89].

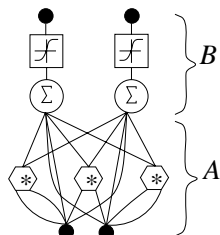


Figure 1: A two-layer high order perceptron.

A high order connection connects a set of neurons in one layer to neurons in the next layer (marked *A* in figure 1). Each connection applies its specific *splicing function* to the activation values of the lower layer. The number of activation values combined by the splicing functions determines the order of the connection and the connection with the highest order determines the order  $\omega$  of the network. A network is called a *full (n-th order) network*, if all possible interlayer connections up to this order are present. The network shown in figure 1 is, for example, a full two layer second order network.

The results of the splicing functions are fed, together with the activation values of the lower layer, into the next layer of neurons (marked *B*). Each neuron consists of a summation unit and an activation function (depicted by a  $\Sigma$  and a symbolized function, respectively).

High order perceptrons can be trained using the backpropagation algorithm, with possible extensions such as a momentum term [Rumelhart-86] or flat spot elimination [Fahlman-88].

From now on in this publication, only two layer high order perceptrons are considered. The splicing function used in this study is multiplication, but other functions are also conceivable.

## 4 The Simulations

A simulation consists of initializing a neural network and applying the online backpropagation algorithm, alternated by convergence tests, until (non-)convergence is observed. A number of simulations starting with the same initial conditions is called an experiment.

The experiments are performed with two major aims: firstly, to see whether the performance of a network changes for different types of initial weight distributions, and secondly to find the optimal initial weight variance, depending on the activation function of the output neurons and the network order.

Each experiment consists of at least 50 simulations. The number of simulations per experiment was increased until the size of the 95% confidence interval for the mean convergence time permitted a sound conclusion. The confidence intervals were calculated under the assumption that the mean convergence time is student-t distributed.

For the simulations performed, a suboptimal learning rate is used, as it is too laborious and computing time consuming to find the optimal learning rate for each combination of data set and network, and as the learning rate and initial weight variance seem to have an independent influence on the learning speed. Because of the suboptimal learning speed, the results do not necessarily allow a comparison between different activation functions and experiments reported elsewhere, as the maximal possible learning rate may differ largely from the one actually used. For example, a third order network has, for the solar data with a shifted/scaled logistic output function, an optimal learning rate of about 0.05. In contrast, the same network and data set, except for using now a standard logistic output function, has an optimal learning rate of about 4.0, and converges in about the same number of iterations.

Some of the convergence criteria chosen in these simulations are rather crude and not related to the task to be solved. This is done in the aim to reduce the high computational expense, which still was several months of Sparc 10 CPU time.

### 4.1 The Data Sets

Most of the data sets used, and shortly described below, are obtained (if not stated otherwise) from an anonymous-ftp server at the University of California [Murphy-94], which also contains further references and documentation. In the following list of data sets, the two numbers in brackets behind the name of the data set are the number of input and output values, respectively.

**Solar (12,1)** contains the sun spot activity for the years 1700 to 1990. The task is to predict the sun spot activity for one of those years, given the activity of the preceding twelve years ( $\hat{=}$  12 real valued inputs). The data are scaled to the interval  $[0, 1]$ .

**Wine (13,3)** is the result of a chemical analysis of wines grown in a region in Italy derived from three different cultivars. The analysis determined the quantities of 13 constituents found in each of the three types of wines. A wine has to be classified using these values, which are scaled to the interval  $[0, 1]$ . The output patterns use boolean values, encoded as +1 and -1.

**CES (2,1)** is the output of the constant elasticity of a substitution production function for thirty pairs of labor and capital input (see [Judge-85], pages 195 and 210). The patterns have two real valued inputs and one real valued output, none of them scaled.

**Servo (12,1)** was created by Karl Ulrich (MIT) in 1986 and contains a very non-linear phenomenon: predicting the rise time of a servomechanism in terms of two (continuous) gain settings and two (discrete) choices of mechanical linkages. The input is coded into two groups of five boolean values each, and two discrete inputs, one assuming four, the other five values. The output is real valued, and like all real valued inputs, scaled to the interval  $[0, 1]$ .

**Vowels (20,5)** is a subset of 300 patterns of the vowels data set, obtainable via ftp from cochlea.hut.fu (130.233.168.48) with the LVQ-package (lvq\_pak). An input pattern consists of 20 unscaled cepstral coefficients obtained from continuous Finnish speech. The task is to determine whether the pronounced phoneme is a vowel, and, in the case it is, which of the five possible ones. The boolean output values are encoded as +1 and -1.

**Auto-mpg (7,1)** concerns city-cycle fuel consumption of cars in miles per gallon, to be predicted in terms of 3 multi-valued discrete and 4 continuous attributes. All values are scaled to the interval  $[0, 1]$  (incomplete patterns have been removed).

**Glass (9,1)** consists of 8 scaled weight percentages of certain oxides and a 7 valued code for the type of glass (window glass, head lamps etc.). The output is the refractive index of the glass, scaled to  $[0, 1]$ .

**Digits (256,10)** consists of 500 handwritten digits (50 patterns for each of the ten digits) of the *NIST Special Database 3* [Garris-92]. Each digit was scaled to fit into an image of 16x16 points, and each pixel is represented by an eight bit value. The input values are scaled to the interval  $[-1, 1]$  and the boolean output values are encoded as +1 and -1.

## 4.2 The Experiments for Multilayer Perceptrons

In the aim to validate and compare the performance of the random weight initialization techniques for multilayer perceptrons mentioned in section 2, a large number of experiments has been performed using the data sets listed in the previous section. The network topology used has one hidden layer which is fully interlayer connected to both input and output layer. The network has no intralayer or supralayer connections, and all activation functions in the hidden and output layer are hyperbolic tangents. No optimization technique was used for training.

For each data set, a sequence of experiments with uniform initial weight distributions of a varying variance were performed (100 simulations per experiment). The outcome of these experiments was used to determine the overall best weight variance as a reference for comparing the random weight initialization schemes of L. Bottou, F. J. Śmieja, G. P. Drago et al., and D. Nguyen et al. It should be noted that D. Nguyen et al. do not seem to use a bias in the output layer. However, neither [Nguyen-90] nor the references mentioned in it state why. To make the simulations fair (leaving out the bias makes learning more difficult), a bias is used in the all simulations reported in this publication.

### 4.2.1 The Results for Multilayer Perceptrons

The outcome of the experiments for the multilayer perceptrons are shown in tables 1 and 2. These tables list in the first column the name of the data set, the number of neurons in the hidden layer  $N_2$ , the convergence criterion  $\epsilon$  (the notation ' $< a$ ' means that the mean square distance between network output and target pattern has to be smaller than  $a$ , and ' $b\%$ ' means that at least  $b$  percent of the patterns must be classified correctly), and the learning rate  $\alpha$ . The subsequent columns labeled with the initial weight variances in table 1 and names in table 2, respectively, contain the outcome of the experiments. The names in table 2 refer to the random weight initialization schemes described in section 2. A single number in these columns corresponds to the mean number of required online learning cycles until convergence (an online learning cycle is a presentation of all patterns with a weight update after each presentation). A number printed in bold face marks the best result in a row and an entry  $p/c$  signifies that the network did not converge in  $p$  percent of the online learning cycles, where a trial is judged as non-convergent if the number of cycles exceeds  $c$ . The rightmost column in table 2 shows the maximal radius  $t_{max}$  of the confidence intervals for the mean number of required learning cycles<sup>3</sup> for the methods listed in this table and a for random weight initialization with a variance of 0.2 (see table 1).

---

<sup>3</sup>The radius of a confidence interval is the difference between the mean and the upper limit of the interval.

	$N_2$	$\epsilon$	$\alpha$	The initial weight variance $\sigma_w^2$													
				$10^{-6}$	0.0001	0.001	0.005	0.01	0.05	0.1	0.2	0.5	0.7	1.0	2.0	3.0	4.0
Solar	5	<0.06	0.3	202	202	195	176	174	165	167	157	148	<b>142</b>	1/500	157	2/500	2/500
Wine	6	90%	0.2	109	105	104	103	104	99.2	100	98.5	97.0	94.3	94.6	<b>92.1</b>	4/500	5/500
CES	4	<0.14	0.1	248	171	131	108	97.9	75.1	63.1	49.7	36.9	<b>35.4</b>	46.9	1/500	1/500	1/500
Servo	3	<0.08	0.1	129	99.0	84.9	74.9	71.0	63.8	64.0	<b>62.8</b>	68.5	73.5	81.3	119	173	4/500
Vowels	20	90%	0.05	158	143	133	125	122	117	113	109	<b>101</b>	102	113	16/800	49/800	53/800
Auto-mpg	3	<0.06	0.3	35.7	33.6	33.9	34.0	33.6	31.9	31.5	<b>31.5</b>	34.7	39.2	44.4	1/500	1/500	1/500
Glass	6	<0.04	0.6	20.9	17.2	15.6	14.1	13.7	12.2	<b>12.0</b>	12.5	16.7	20.0	27.2	48.1	82.8	120
Digits	30	95%	0.3	12.7	12.1	11.6	11.4	<b>11.2</b>	11.6	12.8	15.1	32/100	46/100	25/100	98/100	100/100	100/100

Table 1: Performance of multilayer perceptrons for a uniform weight distribution over the interval  $[-a, a]$ , with  $a = \sqrt{3\sigma^2}$ .

data	$N_2$	$\epsilon$	$\alpha$	Bottou	Wessels	Śmieja	Drago	Nguyen	$t_{max}$
Solar	5	<0.06	0.3	152	<b>146</b>	151	153	162	3.1
Wine	6	90%	0.2	98.0	96.9	<b>96.0</b>	98.4	99.3	1.3
CES	4	<0.14	0.1	40.4	<b>32.5</b>	40.8	44.1	42.7	3.2
Servo	3	<0.08	0.1	56.6	<b>55.6</b>	62.2	59.9	63.2	1.4
Vowels	20	90%	0.05	110	<b>105</b>	111	111	115	1.0
Auto-mpg	3	<0.06	0.3	31.8	<b>30.9</b>	31.4	33.1	31.4	1.2
Glass	6	<0.04	0.6	<b>11.0</b>	11.9	11.7	12.1	11.9	0.4
Digits	30	95%	0.1	<b>11.3</b>	11.7	12.8	12.9	11.4	0.2

Table 2: Random weight initialization with the methods of other authors

#### 4.2.2 Analysis of the simulations for Multilayer Perceptrons

The average convergence behavior of a multilayer perceptron is depicted in figure 2, where region  $A$  indicates the optimum initial weight variances that have been encountered and region  $B$  non-convergence. As the curve is flatter on the left side of the optimal initial weight variance than on the right, the loss in performance is much more tolerable for initial weight variances smaller than the optimal value as compared to bigger variances. Moreover, non-convergence was only encountered for simulations using initial weight variances bigger than the optimal value.

The rather small differences obtained for an initial weight variance of 0.2 as compared to the optimal result, suggests to use this value for a simple weight initialization method. A comparison between the results for this simple method and table 2 shows that the weight initialization method of L. F. A. Wessels et al., which uses the same weight variances as the method of J. W. Boers, performs the best.

Remark: Some of the weight initialization methods presented in section 2 scale the upper and lower bound of the initial random weight interval by the reciprocal square root of the fan-in. This corresponds to scaling the initial weight variance by the reciprocal of the fan-in. Hence, these methods assume a negative correlation between the fan-in of a neuron and the best initial weight variance. This correlation can not be confirmed or rejected by the results of the experiments; more experiments with other data sets are necessary for this.

### 4.3 The Experiments with High Order Perceptrons

The networks used in the simulations are usually full and the biases are initialized with a random value of the same distribution as the weights. The only exception is the network trained on the digits data set. This network includes all first order connections and only second order connections with both inputs corresponding to different pixels in the same row or the same column in the image. This configuration should allow the extraction of sufficient features to learn the digits. Training sessions on the in section 4.1 described digits data set gave an acceptable recognition of untrained digits, despite the small training

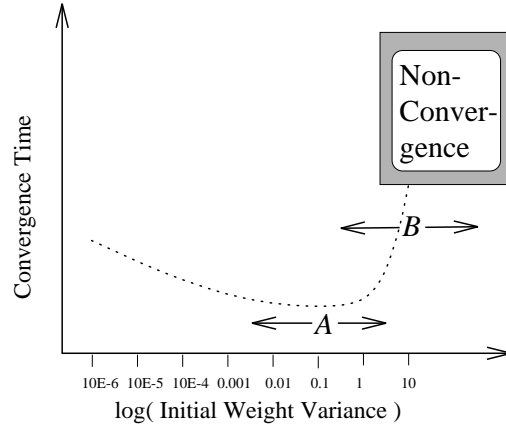


Figure 2: The average behavior of a multilayer perceptron in convergence speed for changing the initial weight variance.

set used.

The three different initial random weight distributions used are: uniform on the interval  $[-a, a]$  (with  $a = \sqrt{3\sigma^2}$ ), normal (restricted to an absolute value of  $3\sigma^2$ ), and uniform over the intervals  $[-2a, -a]$  and  $[a, 2a]$  (with  $a = 3\sigma^2/7$ ) while zero everywhere else. The three types of activation functions used are: a linear  $f_l$ , a hyperbolic tangent  $f_t$ , and a scaled/shifted hyperbolic tangent  $f_{st}$ , shown in figure 3. The use of the function  $f_{st}$  was motivated by several ideas: the scaling in the direction of the y-axis prevents the weights from becoming very big and thus cause the same effect as for example scaling the output data to  $[-0.9, 0.9]$  and the change of the steepness and the shifting of the sigmoid in the direction of the x-axis were used to force the outcome of the summation step in the neurons to be in the interval  $[0, 1]$ . Also, experiments with this activation function where performed to see, whether a relation between a deformation of the activation function and the optimal initial weight range exists. The only optimization technique applied to speed up learning is flat-spot elimination.

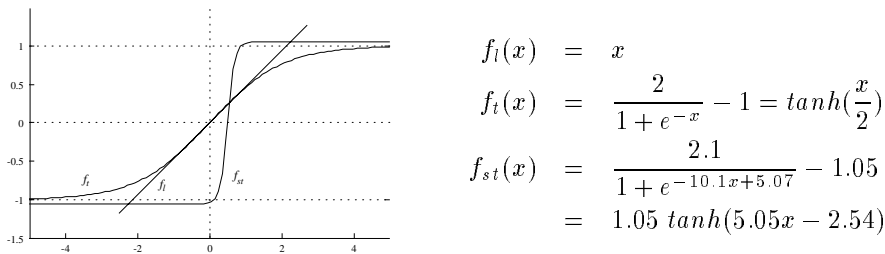


Figure 3: The activation functions

#### 4.3.1 The Results of the Simulations for High Order Perceptrons

Tables 3, 4, and 5 show, besides entries already explained in section 4.2.1, the order  $\omega$  of the fully interlayer connected network and the activation function  $f$ .

#### 4.3.2 Analysis of the Simulations for High Order Perceptrons

The minimal convergence times for all three initial weight distributions show no difference of statistical significance. The average behavior of the learning time as a function of the initial weight variance, which is depicted in figure 4, is explained as follows. The main difference for the three distributions is the value of the weight variance where the convergence time starts increasing drastically. This “edge” (point A) is roughly at the same location for both the uniform distribution and the uniform distribution



				The initial weight variance $\sigma^2$												
$\omega$	$f$	$\epsilon$	$\alpha$	$10^{-6}$	0.0001	0.001	0.005	0.01	0.05	0.1	0.2	0.5	0.7	1.0	conf.	
Solar	1	$f_l$	<0.08	0.1	<b>1.3</b>	1.4	<b>1.3</b>	1.4	1.5	1.9	2.5	3.7	6.1	8.6	10.2	1.1, 1.4
	2	$f_l$	<0.06	0.1	5.6	<b>5.4</b>	5.6	6.1	6.6	10.9	15.6	23.4	46.5	58.3	82.0	5.2, 5.7
	3	$f_l$	<0.06	0.1	<b>4.5</b>	<b>4.5</b>	4.9	4.9	5.6	9.4	15.4	24.6	53.8	81.9	115.1	4.3, 4.7
	1	$f_t$	<0.08	0.7	<b>4.5</b>	<b>4.5</b>	<b>4.5</b>	<b>4.5</b>	<b>4.5</b>	4.7	4.7	4.9	5.3	5.8	6.2	4.4, 4.6
	2	$f_t$	<0.06	0.7	37.5	<b>37.1</b>	37.2	37.4	37.6	38.2	39.0	39.9	45.3	48.4	54.9	36.9, 37.3
	3	$f_t$	<0.06	0.7	<b>22.7</b>	<b>22.7</b>	22.8	22.8	22.8	23.5	24.4	26.2	31.5	33.8	40.5	22.5, 22.8
Wine	1	$f_t$	80%	0.5	242	242	242	242	241	241	241	240	237	236	235, 237	
	2	$f_t$	90%	0.5	213	215	214	215	214	214	214	210	<b>200</b>	1/500	1/500	197, 204
	1	$f_{st}$	80%	0.02	<b>32.4</b>	32.6	32.6	32.7	33.5	8/300	25/300	41/300	68/300	76/300	89/300	32.2, 32.6
	2	$f_{st}$	90%	0.02	<b>24.1</b>	<b>24.1</b>	<b>24.1</b>	24.2	25.6	14/300	34/300	57/300	77/300	70/300	78/300	23.8, 24.4
CES	1	$f_l$	<0.14	0.1	6.0	6.0	6.2	6.0	6.0	<b>5.9</b>	6.4	6.4	8.4	9.3	10.4	5.6, 6.1
	2	$f_l$	<0.08	0.1	57.6	<b>57.3</b>	57.6	58.5	58.0	60.9	63.7	66.0	81.7	81.0	91.2	56.7, 57.9
	1	$f_t$	<0.14	0.2	12.7	12.7	12.7	12.7	12.6	12.5	12.7	<b>12.4</b>	12.7	13.4	14.5	12.0, 12.6
	2	$f_t$	<0.08	0.2	<b>158</b>	159	<b>158</b>	<b>158</b>	159	159	160	163	167	173	4/300	157, 159
	1	$f_{st}$	<0.14	0.01	<b>14.1</b>	14.6	14.7	14.6	15.1	3/300	10/300	13/300	23/300	27/300	37/300	13.7, 14.6
	2	$f_{st}$	<0.08	0.01	23.5	<b>23.1</b>	23.6	26.2	28.6	5/300	11/300	19/300	44/300	51/300	73/300	22.7, 23.5
Servo	1	$f_l$	<0.14	0.01	2.1	<b>2.0</b>	2.1	2.7	3.44	5.8	8.0	9.7	11.3	15.6	16.0	1.9, 2.2
	2	$f_l$	<0.08	0.01	62	<b>61</b>	<b>61</b>	62	64	71	83	99	159	186	209	61, 62
	1	$f_t$	<0.14	0.03	19.6	19.6	19.6	19.4	19.2	<b>18.5</b>	18.7	19.7	23.7	27.4	35.2	18.1, 19.0
	2	$f_t$	<0.08	0.03	<b>168</b>	<b>168</b>	<b>168</b>	<b>168</b>	<b>168</b>	171	172	183	217	231	290	168, 168
Vowels	2	$f_l$	65%	0.03	<b>30.6</b>	31.7	31.5	32.4	32.3	36.3	42.7	51.9	73.2	82.6	103.3	29.8, 31.4
	1	$f_t$	80%	0.5	64.4	64.5	64.9	64.5	64.6	64.6	63.9	63.8	<b>63.7</b>	63.9	64.3	63.3, 64.1
	2	$f_t$	90%	0.5	15.3	15.3	15.4	15.3	15.3	15.2	15.3	<b>15.1</b>	20.0	25.2	4/400	14.9, 15.3
	1	$f_{st}$	80%	0.03	9.0	<b>8.9</b>	9.0	9.0	9.8	39.1	6/400	74/400	86/400	98/400	88/400	8.6, 9.1
	2	$f_{st}$	90%	0.03	<b>31.4</b>	42.4	32.6	2/1000	5/1000	82/1000	92/1000	96/1000	98/1000	98/1000	100/1000	27.0, 40.6
Auto-mpg	2	$f_l$	<0.06	0.03	10.9	<b>10.8</b>	10.9	11.4	11.9	16.3	24.2	32.9	75.3	3/200	24/200	10.7, 11.0
	2	$f_t$	<0.06	0.1	<b>22.9</b>	23.0	<b>22.9</b>	23.7	23.2	24.5	26.3	28.9	38.8	45.6	66.8	22.5, 23.3
Glass	1	$f_l$	<0.04	0.1	4.1	<b>4.0</b>	<b>4.0</b>	<b>4.0</b>	<b>4.0</b>	4.4	5.0	5.8	7.6	8.2	9.3	3.8, 4.1
	2	$f_l$	<0.03	0.005	<b>13.3</b>	<b>13.3</b>	15.5	30.7	54.6	221	52/400	100/400	100/400	100/400	100/400	13.1, 13.5
Digits	2r	$f_l$	95%	0.0002	<b>44.7</b>	46.8	63.0	119	172	100/100	100/100	100/100	100/100	100/100	100/100	44.1, 45.3

				$\sigma^2$			
$\omega$	$f$	$\epsilon$	$\alpha$	2.0	5.0	10.0	
Wine	1	$f_t$	80%	0.5	<b>234</b>	4/500	20/500
	2	$f_t$	90%	0.5	4/500	24/500	58/500

				$\sigma^2$		
$\omega$	$f$	$\epsilon$	$\alpha$	$10^{-10}$	$10^{-8}$	
Digits	2r	$f_l$	95%	0.0002	45.2	45.2

Table 3: Performance of high order perceptrons for a uniform weight distribution over the interval  $[-a, a]$ , with  $a = \sqrt{3\sigma^2}$ .

					The initial weight variance $\sigma^2$											
$\omega$	$f$	$\epsilon$	$\alpha$	$10^{-6}$	0.0001	0.001	0.005	0.01	0.05	0.1	0.2	0.5	0.7	1.0	conf.	
Solar	1	$f_l$	<0.08	0.1	<b>1.3</b>	<b>1.3</b>	1.5	<b>1.3</b>	1.5	1.5	1.6	2.5	5.6	7.4	10.2	1.2, 1.4
	2	$f_l$	<0.06	0.1	5.5	<b>5.4</b>	5.5	5.5	5.5	6.3	7.9	14.8	39.0	55.4	79.8	5.3, 5.5
	3	$f_l$	<0.06	0.1	<b>4.5</b>	<b>4.5</b>	<b>4.5</b>	4.8	4.8	5.6	8.4	14.4	45.1	68.4	111.1	4.4, 4.7
	1	$f_t$	<0.08	0.7	4.5	4.5	4.5	4.5	<b>4.4</b>	4.5	4.6	4.7	5.2	5.8	6.4	4.4, 4.5
	2	$f_t$	<0.06	0.7	37.3	37.4	37.4	<b>37.2</b>	<b>37.2</b>	37.3	37.7	38.9	43.9	48.1	54.6	37.3, 37.6
	3	$f_t$	<0.06	0.7	<b>22.6</b>	<b>22.6</b>	22.7	22.7	22.7	22.9	23.0	24.7	29.6	34.2	38.8	22.5, 22.8
Wine	1	$f_t$	80%	0.5	242	242	242	242	242	242	242	241	238	240	236	227, 239
	2	$f_t$	90%	0.5	215	215	215	213	214	213	216	215	204	203	<b>199</b>	196, 202
	1	$f_{st}$	80%	0.02	32.5	32.4	<b>32.3</b>	32.5	32.5	32.9	1/300	24/300	67/300	72/300	79/300	32.0, 32.6
	2	$f_{st}$	90%	0.02	<b>23.9</b>	<b>23.9</b>	24.3	24.0	24.1	24.3	1/300	25/300	76/300	75/300	83/300	23.5, 24.2
CES	1	$f_l$	<0.14	0.1	<b>6.0</b>	<b>6.0</b>	<b>6.0</b>	<b>6.0</b>	6.1	<b>6.0</b>	<b>6.0</b>	6.4	7.5	8.7	9.7	5.9, 6.1
	2	$f_l$	<0.08	0.1	57.7	57.8	57.7	57.8	<b>57.6</b>	58.1	59.5	64.2	74.4	81.6	98.7	57.4, 57.8
	1	$f_t$	<0.14	0.2	12.4	12.9	12.5	12.7	12.8	12.8	12.5	12.5	<b>12.1</b>	13.4	14.7	11.4, 12.8
	2	$f_t$	<0.08	0.2	158	158	158	158	158	<b>157</b>	160	160	167	2/300	6/300	155, 158
	1	$f_{st}$	<0.14	0.01	15.0	14.8	14.5	<b>13.9</b>	14.7	14.6	21.7	10/300	20/300	51/300	50/300	13.4, 14.3
	2	$f_{st}$	<0.08	0.01	23.4	23.2	23.1	23.1	<b>23.0</b>	27.3	43.7	8/300	42/300	58/300	62/300	22.6, 23.4
servo	1	$f_l$	<0.14	0.01	2.1	<b>1.9</b>	2.0	2.1	2.1	2.8	4.7	7.7	12.7	13.3	14.7	1.8, 2.0
	2	$f_l$	<0.08	0.01	<b>62</b>	<b>62</b>	<b>62</b>	<b>62</b>	<b>62</b>	63	66	81	145	164	231	61, 62
	1	$f_t$	<0.14	0.03	19.6	19.6	19.7	19.6	19.6	<b>19.3</b>	<b>19.3</b>	19.5	22.5	25.6	34.2	19.1, 19.4
	2	$f_t$	<0.08	0.03	168	168	168	168	168	<b>167</b>	170	173	211	237	2/500	166, 168
Vowels	2	$f_l$	65%	0.03	31.7	32.4	32.1	<b>30.0</b>	31.5	32.4	34.2	41.7	68.1	80.9	100.8	28.9, 31.1
	1	$f_t$	80%	0.5	64.4	64.6	64.5	64.7	65.0	64.4	64.9	64.0	63.7	<b>63.3</b>	63.7	62.6, 64.0
	2	$f_t$	90%	0.5	15.7	15.4	15.5	15.3	15.4	15.5	<b>15.2</b>	15.4	16.2	24.7	4/400	15.0, 15.4
	1	$f_{st}$	80%	0.03	<b>8.6</b>	8.8	8.9	8.8	8.9	9.3	15.4	8/400	86/400	95/400	95/400	8.4, 8.8
	2	$f_{st}$	90%	0.03	35.0	<b>31.4</b>	35.9	2/1000	2/1000	2/1000	26/1000	98/1000	98/1000	98/1000	98/1000	24.7, 38.0
Auto-mpg	2	$f_l$	<0.06	0.03	10.8	10.9	10.8	<b>10.7</b>	<b>10.7</b>	11.3	13.6	23.0	69.1	3/200	28/200	10.6, 10.9
	2	$f_t$	<0.06	0.1	22.8	<b>22.7</b>	22.9	22.9	22.8	23.1	24.0	26.6	37.0	45.0	67.3	22.5, 22.9
Glass	1	$f_l$	<0.04	0.1	4.2	4.1	<b>4.0</b>	4.1	<b>4.0</b>	4.3	4.2	4.9	6.9	8.1	9.1	3.8, 4.1
	2	$f_l$	<0.03	0.005	<b>13.3</b>	<b>13.3</b>	<b>13.3</b>	13.4	13.8	39.4	128	48/400	100/400	100/400	100/400	13.2, 13.4

					$\sigma^2$			
$\omega$	$f$	$\epsilon$	$\alpha$		2.0	5.0	10.0	
Wine	1	$f_t$	80%	0.5	<b>234</b>	1/500	16/500	
	2	$f_t$	90%	0.5	4/500	16/500	36/500	

Table 4: Performance of high order perceptrons for a normal distribution restricted to the interval  $[-3\sigma^2, 3\sigma^2]$ .

				$\sigma^2$												conf.
$\omega$	$f$	$\epsilon$	$\alpha$	$10^{-6}$	0.0001	0.001	0.005	0.01	0.05	0.1	0.2	0.5	0.7	1.0		
Solar	1	$f_l$	<0.08	0.1	<b>1.2</b>	<b>1.2</b>	1.4	1.4	1.4	2.1	2.7	3.5	6.5	9.3	11.5	1.1, 1.3
	2	$f_l$	<0.06	0.1	5.4	<b>5.3</b>	5.6	5.8	6.7	10.1	15.0	24.1	46.0	60.2	83.5	5.2, 5.5
	3	$f_l$	<0.06	0.1	<b>4.5</b>	4.6	4.7	5.2	5.8	10.2	14.7	24.4	56.2	77.8	114	4.3, 4.6
	1	$f_t$	<0.08	0.7	4.5	<b>4.4</b>	4.5	4.5	4.6	4.8	4.7	4.8	5.4	5.5	6.5	4.3, 4.5
	2	$f_t$	<0.06	0.7	37.4	<b>37.1</b>	<b>37.1</b>	37.4	37.3	37.9	39.1	41.0	47.9	48.6	53.7	36.9, 37.4
	3	$f_t$	<0.06	0.7	22.7	<b>22.6</b>	<b>22.8</b>	<b>22.6</b>	23.0	23.5	24.4	26.3	31.3	34.8	38.9	22.5, 22.8
Wine	1	$f_t$	80%	0.5	242	242	242	242	241	241	241	240	238	237	231, 234	
	2	$f_t$	90%	0.5	215	215	215	214	214	214	210	208	<b>198</b>	<b>198</b>	194, 202	
	1	$f_{st}$	80%	0.02	<b>32.6</b>	<b>32.6</b>	<b>32.6</b>	32.8	33.3	9/300	16/300	51/300	72/300	76/300	83/300	32.3, 32.9
	2	$f_{st}$	90%	0.02	24.3	<b>23.9</b>	23.8	<b>23.9</b>	25.8	15/300	40/300	58/300	69/300	64/300	79/300	23.6, 24.3
CES	1	$f_l$	<0.14	0.1	6.0	6.0	<b>5.9</b>	6.1	6.0	6.0	6.6	9.0	9.8	10.5	5.7, 6.0	
	2	$f_l$	<0.08	0.1	<b>56.7</b>	58.0	57.7	57.8	58.5	59.8	65.8	70.5	85.1	84.1	94.7	56.1, 57.3
	1	$f_t$	<0.14	0.2	12.6	12.7	12.8	12.5	12.5	12.6	<b>12.3</b>	12.5	12.8	12.8	14.5	11.8, 12.7
	2	$f_t$	<0.08	0.2	159	<b>158</b>	<b>158</b>	<b>158</b>	159	160	162	164	176	173	7/300	157, 159
	1	$f_{st}$	<0.14	0.01	14.6	<b>14.0</b>	14.8	15.0	15.6	37.8	8/300	13/300	25/300	25/300	31/300	13.5, 14.5
	2	$f_{st}$	<0.08	0.01	23.4	<b>22.9</b>	23.7	26.1	30.1	3/300	11/300	29/300	51/300	65/300	74/300	22.3, 23.4
Servo	1	$f_l$	<0.14	0.01	<b>2.1</b>	2.3	<b>2.1</b>	2.8	3.0	5.9	7.7	10.0	12.7	13.5	17.3	1.9, 2.3
	2	$f_l$	<0.08	0.01	61.4	61.6	<b>61.1</b>	62.6	64.9	69.3	83.5	100	153	195	225	60.4, 61.8
	1	$f_t$	<0.14	0.03	19.6	19.6	19.5	19.6	19.5	<b>19.1</b>	19.2	19.7	23.4	27.7	33.7	18.7, 19.4
	2	$f_t$	<0.08	0.03	168	168	<b>167</b>	168	168	172	176	182	221	253	309	167, 168
Vowels	2	$f_l$	65%	0.03	30.9	32.9	<b>30.8</b>	32.9	33.4	38.4	41.1	50.8	73.0	85.9	101.2	29.5, 32.2
	1	$f_t$	80%	0.5	64.3	64.9	64.4	63.6	64.2	63.9	63.7	64.0	63.4	<b>63.2</b>	63.7	62.5, 63.9
	2	$f_t$	90%	0.5	15.3	15.3	15.2	15.4	15.6	<b>14.9</b>	15.5	15.6	17.9	21.9	5/400	14.6, 15.1
	1	$f_{st}$	80%	0.03	9.0	<b>8.7</b>	8.9	9.3	9.8	37.2	8/400	60/400	88/400	98/400	96/400	8.5, 9.0
2	$f_{st}$	90%	0.03	<b>28.6</b>	43.6	40.3	1/1000	2/1000	76/1000	95/1000	99/1000	98/1000	98/1000	100/1000	23.3, 33.9	
Auto-mpg	2	$f_l$	<0.06	0.03	<b>10.8</b>	<b>10.8</b>	10.9	11.5	11.8	17.0	22.5	34.0	74.0	10/200	28/200	10.5, 11.1
	2	$f_t$	<0.06	0.1	23.0	22.8	<b>22.7</b>	23.4	23.5	24.7	26.4	28.2	43.5	49.1	69.8	22.3, 23.1
Glass	1	$f_l$	<0.04	0.1	4.0	<b>3.9</b>	4.1	4.2	4.1	4.5	5.0	6.0	7.2	8.7	9.6	3.7, 4.1
	2	$f_l$	<0.03	0.005	<b>13.4</b>	<b>13.4</b>	15.3	29.5	52.2	223	56/400	100/400	100/400	100/400	100/400	13.2, 13.6

				$\sigma^2$			
$\omega$	$f$	$\epsilon$	$\alpha$	2.0	5.0	10.0	
Wine	1	$f_t$	80%	0.5	<b>233</b>	4/500	32/500
	2	$f_t$	90%	0.5	4/500	12/500	34/500

Table 5: Performance of high order perceptrons for a uniform weight distribution over the intervals  $[-2a, -a]$  and  $[a, 2a]$ , with  $a = \sqrt{3\sigma^2/7}$ .

over two intervals, but slightly shifted to higher variances for the normal distribution. As the optimal weight variance (point  $B$ ) is similarly displaced, the performance of two networks, initialized with two different weight distributions of the same variance, is difficult to compare. This might explain the better performance for a Gaussian initial weight distribution in the report of P. Haffner. For the various combinations of data set, network order, etc., the optimal weight variance was encountered in region  $D$ , whereas non-convergence was, if at all, only observed in region  $C$ .

As the three different initial weight distributions yield no significant difference in network performance, only the commonly used uniform distribution is considered from now on. For the shifted/scaled logistic and the linear activation functions, the best fixed weight variance is about  $10^{-4}$  (which corresponds to an interval  $[-0.017, 0.017]$ ). For the logistic activation function, the best value for the weight variance depends a lot on data set and network order. In general, the performance with optimal initial weight variance differs not much more than about 10% from the results obtained with a variance of  $10^{-4}$  or even smaller. Therefore a variance of  $10^{-4}$  may be used as a simple application independent random weight initialization scheme. This initialization scheme is also justified by a smaller risk: a network performs nearly as good for an initial weight variance smaller than the optimum. The loss in performance for choosing the initial weight variance too small is much less significant than it is for multilayer perceptrons.

The experiments confirm also that the data set itself has a large influence on the optimal initial weight variance: for the solar, wine, and servo data sets, the networks have about the same size for the same order, but the optimal value for the weight variance differs a lot for the network with the logistic

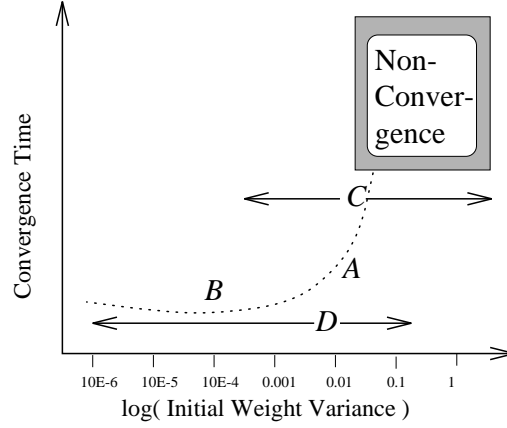


Figure 4: Average behavior of a higher order perceptron in convergence time for changing initial weight variance

activation function. Further, the optimal value for the initial weights remained for some data sets nearly unchanged for different net orders or even different activation functions, while it changes greatly for other sets. It remains the question, which attribute of the data sets causes this behavior.

## 5 Weight Initialization Techniques for Three Layer Perceptrons Applied to High Order Perceptrons

The most remarkable fact is that the rules based on observations (the “rules of thumb”), which seem to perform well for multilayer perceptrons, do not apply directly to high order perceptrons. Uniform weight distributions with an initial weight range of  $[-1, 1]$  (which corresponds to a variance of about 0.33) or bigger are definitely a poor choice for most of the examples considered. The mean convergence time is for some of the examples more than four times higher than the best initial weight range found or, even worse, they do not converge in a reasonable time.

The approaches using an interval  $[-a/\sqrt{d_{in}}, a/\sqrt{d_{in}}]$  with a more or less arbitrary constant  $a$ , do not outperform the weight initialization with a fixed variance of 0.0001 or vice versa. Nevertheless, one would expect that more sophisticated methods for random weight initialization perform better than a scheme with fixed initial weight variance.

One such a sophisticated scheme was described by L. F. A. Wessels [Wessels-92], and is here recalculated for a second order net with linear activation functions. The method tries to initialize the weights in such a way that the variance  $\sigma_y^2$  of the network output is equal to the expected variance  $\sigma_y^2$  of the target patterns.

A random weight initialization allows the calculation of the variance  $\sigma_y^2$  depending on the network topology ( $N_1$  is the number of neurons in the input layer,  $w_i$  a weight,  $x_i$  an input value):

$$\begin{aligned}
 \sigma_y^2 &= E[\hat{y}^2] - E^2[\hat{y}] = E[\hat{y}^2] \quad \text{as the weights are independent from the input values} \\
 &= E\left[\sum_{i,j=1}^{N_1} w_i w_j x_i x_j + \sum_{i,j,k,l=1}^{N_1} w_{ij} w_{kl} x_i x_j x_k x_l\right] \\
 &= \sum_{i=1}^{N_1} E[w_i^2] E[x_i^2] + \sum_{i,j=1}^{N_1} E[w_{ij}^2] E[x_i^2] E[x_j^2] \quad \text{as } E[w_{ij} w_{jk}] = 0 \text{ if } i \neq k \text{ or } j \neq l \\
 &= E[w^2] \left( \sum_{i=1}^{N_1} E[x_i^2] + \sum_{i,j=1; i \neq j}^{N_1} E[x_i^2] E[x_j^2] + \sum_{i=1}^{N_1} E[x_i^4] \right) \quad \text{all weights have the same distribution}
 \end{aligned}$$

$$\begin{aligned}
&= E[w^2] \left( \sum_{i=1}^{N_1} \frac{1}{12} + \sum_{i,j=1; i \neq j}^{N_1} \left(\frac{1}{12}\right)^2 + \sum_{i=1}^{N_1} \frac{1}{80} \right) \\
&= \frac{5N_1^2 + 133N_1}{720} E[w^2]
\end{aligned}$$

This assumes that all  $w_i$  and all  $x_i$  are independent random variables in the interval  $[-1, 1]$  and  $[0, 1]$  respectively, and therefore  $E[x^2] = 1/12$  and  $E[x^4] = 1/80$ , as shown by the following equations:

$$\begin{aligned}
E[x^2] &= \int_{-\infty}^{\infty} (x - \frac{1}{2})^2 F_{uniform [0,1]} dx = \int_0^1 (x - \frac{1}{2})^2 dx = \frac{1}{12} \\
E[x^4] &= \int_{-\infty}^{\infty} (x - \frac{1}{2})^4 F_{uniform [0,1]} dx = \int_0^1 (x - \frac{1}{2})^4 dx = \frac{1}{80}
\end{aligned}$$

The variance  $\sigma_w^2$  of a variable  $w$  is equal to the expectation of its squared value  $E[w^2]$ :

$$E[w^2] = \frac{\int_{-\infty}^{\infty} w^2 * F_{uniform, \sigma}(w) dw}{\int_{-\infty}^{\infty} F_{uniform, \sigma}(w) dw} = \frac{\int_{-\sqrt{3\sigma_w^2}}^{\sqrt{3\sigma_w^2}} w^2 dw}{\int_{-\sqrt{3\sigma_w^2}}^{\sqrt{3\sigma_w^2}} 1 dw} = \frac{\frac{1}{3} \cdot 2 \cdot \sqrt{3\sigma_w^2}^3}{2 \cdot \sqrt{3\sigma_w^2}} = \sigma_w^2$$

Therefore the initial condition, equation  $\sigma_y^2 = \sigma_w^2$ , can be fulfilled by setting

$$\sigma_w^2 = \frac{720}{5N_1^2 + 133N_1} \sigma_y^2$$

Where  $\sigma_y^2$  is equal to  $1/3$  for boolean data and data which are scaled to the interval  $[-1, 1]$ , and  $1/12$  for data which are scaled to the interval  $[0, 1]$ .

Similar formulas can easily be calculated for networks of different orders. The evaluation of these formulas gives an optimal initial weight variance between 0.3 and 0.02 for the examples considered in this publication. These values are usually much too high as compared to the results listed in the tables. No effort was therefore done to solve the same problems for the other (non-linear) activation functions.

For the activation functions  $f_t$  and  $f_{st}$  the heuristic of L. Bottou comes the closest to the optimal weight variance and may therefore be considered as being better. Comparing the variances suggested by his heuristic to table 3, one finds that his method may be improved by using about one tenth of the suggested variance.

## 6 Conclusion

The experiments show that a suitable and convenient weight initialization method for high order perceptrons<sup>4</sup> with identity activation function<sup>5</sup> is a random initialization with a rather small variance of about  $10^{-4}$  (which corresponds to a weight range of  $[-0.017, 0.017]$ ). If a hyperbolic tangent is used as activation function, the best performance is obtained for the interval  $[-a/\sqrt{d_{in}}, a/\sqrt{d_{in}}]$ , where  $a$  is chosen such that the weight variance corresponds to one third of the distance between the points of the maximal curvature of the activation function<sup>6</sup> (this is approximately 0.8 for an unscaled hyperbolic tangent with steepness one). The ‘‘rules of thumb’’ which perform well for multilayer perceptrons are not suitable for high order perceptrons (which can be explained by their different topologies).

The steepness (and/or the horizontal shift) of the activation function has a big influence on the convergence time of high order perceptrons. A mathematical study, partly inspired by this research,

<sup>4</sup>The results of this study apply in part to standard (first order) perceptrons, since high order perceptrons are a generalization of standard perceptrons.

<sup>5</sup>A linear activation function of steepness one.

<sup>6</sup>This results in about one tenth of the weight variance.

showed that this is indeed the case if the steepness of the activation function is changed (assuming the initial weight range is adapted) [Thimm-94.2].

On the other hand, the shape of the initial weight distribution of three rather different distributions showed no or only very little effect on the optimal convergence time of high order perceptrons. The main effect observed is a dislocation of the optimal value for the initial weight variance. There is consequently no preference for one of the three distributions as the optimal learning speeds are similar.

For multilayer perceptrons with one hidden layer, the weight initialization method of L. F. A. Wessels et al. performed on average the best, but the performance of the methods proposed by J. W. Boers et al., L. Bottou, F. J. Śmieja, G. P. Drago et al., and D. Nguyen et al. are nearly as good. Despite the fact that the method of L. F. A. Wessels et al. and L. Bottou apply the same initial weight variances for the experiments performed for this publication, the first should be preferred, as it scales the initial weight variances depending on the activation function (due to the calculation of the network output variance).

The experiments show that the best initial weight variance for both types of neural networks is determined by the data set. Consequently, some reasoning on the data set has to be included in the determination of this value, if better values than those proposed in this publication are desired. On the other hand, an initial weight variance close to the optimal value is often acceptable, as the impact on the number of required learning cycles is not too big for small deviations in variance. In general, the loss in convergence speed for both types of neural networks is bigger when too high a variance is chosen than when too small a variance is chosen, as compared to the optimal value.

The evaluation of the experiments performed in order to find the best weight initialization scheme for high order and multilayer perceptrons includes the calculation of confidence intervals for the mean convergence time. This is a much more reliable measure than simply counting the number of simulations performed, as used in most other publications. Moreover, the simulations showed that some data sets require a bigger amount of simulations for a sufficiently small size of the confidence interval than others. In the experiments performed in this research, these numbers varied between 50 and 2,000. The authors encourage other researchers to report their results in a similar way.

## Acknowledgments

The authors want to thank the Institute for Logic, Complexity, and Deduction Systems at the University of Karlsruhe for the supply of computing power, which allowed the performance of the simulations.

## References

- [Bellido-93] I. Bellido and E. Fiesler. **Do Backpropagation Trained Neural Networks Have Normal Weight Distributions?** In Stan Gielen and Bert Kappen (eds.), *ICANN '93; Proceedings of the International Conference on Artificial Neural Networks*, pp. 772–775, London, U.K., 1993. Springer-Verlag.
- [Boers-92] E. J. W. Boers and H. Kuiper. **Biological Metaphors and the Design of Modular Artificial Neural Networks.** Master's thesis, Leiden University, Leiden, The Netherlands, Aug. 1992.
- [Bottou-88] L.-Y. Bottou. **Reconnaissance de la Parole par Reseaux Multi-Couches.** In *Neuro-Nîmes'88; Proceedings of the International Workshop on Neural Networks and Their Applications*, pp. 197–217, 1988. ISBN: 2-906899-14-3
- [Chen-91] C. L. Chen and R. S. Nutter. **Improving the Training Speed of Three-Layer Feedforward Neural Nets by Optimal Estimation of the Initial Weights.** In *International Joint Conference on Neural Networks*, vol. 3, pp. 2063–2068. IEEE, 1991.
- [Denoeux-93] T. Denoeux and R. Lengellé. **Initializing Back Propagation Networks with Prototypes.** *Neural Networks*, vol. 6, pp. 351–363, Pergamon Press Ltd., 1993.
- [Drago-92] G. P. Drago and S. Ridella. **Statistically Controlled Activation Weight Initialization (SCAWI).** *IEEE Transactions on Neural Networks*, vol. 3, num. 4, pp. 627–631, Jul. 1992.

- [Fahlman-88] S. E. Fahlman. **An Empirical Study of Learning Speed in Backpropagation Networks**. Technical Report CMU-CS-88-162, School of Computer Science, Carnegie Mellon University, Pittsburgh, PA, Sep. 1988.
- [Fiesler-94] E. Fiesler. **Neural Network Classification and Formalization**. In J. Fulcher (ed.), *Computer Standards & Interfaces*, vol. 16, num. 3, special issue on Neural Network Standardization, pp. 231–239. North-Holland/Elsevier, 1994. ISSN: 0920–5489
- [Garris-92] M. D. Garris and R. A. Wilkinson. **NIST Special Database 3**. National Institute of Standards and Technology, Advanced System Division, Image Recognition Group, Feb. 1992.
- [Hertz-91] J. Hertz, A. Krogh, and R. G. Palmer. **Introduction to the Theory of Neural Computation**, vol. I. Addison Wesley, 1991. ISBN: 0-201-51560-1
- [Haffner-88] P. Haffner, A. Waibel, H. Sawai, and K. Shikano. **Fast Back-Propagation Learning Methods for Neural Networks in Speech**. Technical Report TR-1-0058, ATR Interpreting Telephony Research Laboratories, 1988.
- [Judge-85] G. G. Judge, W. E. Griffiths, R. Carter Hill, and T.-C. Lee. **The Theory and Practice of Econometrics**. Wiley Series in Probability and mathematical statistics. John Wiley and Sons, 2nd edition, 1985.
- [Kolen-90] J. F. Kolen and J. B. Pollack. **Back Propagation is Sensitive to Initial Conditions**. Technical Report TR 90-JK-BPSIC. Laboratory for Artificial Intelligence Research, Computer and Information Science Department, 1990.
- [Kim-91] Y. K. Kim and J. B. Ra. **Weight Value Initialization for Improving Training Speed in the Backpropagation Network**. In *International Joint Conference on Neural Networks*, vol. 3, pp. 2396–2401. IEEE, 1991.
- [Lee-86] Y. C. Lee, G. Doolen, H. Chen, G. Sun, T. Maxwell, H. Lee, and C. L. Giles. **Machine Learning Using a Higher Order Correlation Network**. *Physica D: Nonlinear Phenomena*, vol. 22, pp. 276–306, 1986. ISSN: 0167-2789
- [Lee-93] Y. Lee, S.-H. Oh, and M. W. Kim. **An Analysis of Premature Saturation in Back Propagation Learning**. *Neural Networks*, vol. 6, pp. 719–728, 1993.
- [Murphy-94] P. M. Murphy and D. W. Aha (Librarians). **UCI Repository of machine learning databases** [Machine-readable data repository], anonymous-ftp access ics.uci.edu: pub/machine-learning-databases, 1994.
- [Nguyen-90] D. Nguyen and B. Widrow. **Improving the Learning Speed of 2-Layer Neural Networks by Choosing Initial Values of the Adaptive Weights**. In *Proceedings of the International Joint Conference on Neural Networks (IJCNN) San Diego*, vol. III, pp. 21–26, Edward Brothers, 1990.
- [Pao-89] Y.-H. Pao. **Adaptive Pattern Recognition and Neural Networks**. Addison-Wesley Publishing Company, Inc., Reading, Mass., 1989. ISBN: 0-201-12584-6
- [Platt-91] J.C. Platt. **Learning by Combining Memorization and Gradient Descent**. In R. P. Lippman et al. (eds.), *Advances in Neural Information Processing Systems*, vol. III, pp. 714–720. Morgan Kaufmann, San Mateo, 1991.
- [Rumelhart-86] D. E. Rumelhart, J. L. McClelland, and the PDP Research Group. **Parallel Distributed Processing: Explorations in the Microstructure of Cognition**, vol. 1: Foundations. The MIT Press, Cambridge, Mass., 1986. ISBN: 0-262-18120-7
- [Smieja-91] F. J. Śmieja. **Hyperplane “Spin” Dynamics, Network Plasticity and Back-Propagation Learning**. GMD report, GMD, St. Augustin, Germany, Nov. 28, 1991.
- [Thimm-94.1] G. Thimm, R. Grau, and E. Fiesler. **Modular Object-Oriented Neural Network Simulators and Topology Generalizations**. In M. Marinaro and P. G. Morasso (eds.), *Proceedings of the International Conference on Artificial Neural Networks (ICANN 94)*, vol. 1, pp. 747–750, London, U.K., 1994. Springer-Verlag. ISBN: 3-540-19887-3
- [Thimm-94.2] G.. Thimm, P. Moerland, and E. Fiesler. **The Learning Rate and the Gain of the Activation Function in Backpropagation Neural Networks are Exchangeable**. Submitted to Neural Computation. See also P. Moerland, G. Thimm, and E. Fiesler. **Results on the Steepness in Backpropagation Neural Networks**. In Marc Aguilar (ed.), *Proceedings of the '94 SIPAR-Workshop on Parallel and Distributed Computing*, Inst. of Informatics, University Pérolles, Chemin du Musée 3, Fribourg, Switzerland, pp. 91–94, Oct. 1994. SI Group for Parallel Systems.

- [Wessels-92] L. F. A. Wessels and E. Barnard. **Avoiding False Local Minima by Proper Initialization of Connections.** *IEEE Transactions on Neural Networks*, vol. 3, num. 6, pp. 899–905, Nov. 1992.
- [Watrous-93] R. L. Watrous and G. M. Kuhn. **Some Considerations on the Training of Recurrent Neural Networks for Time-Varying Signals.** In M. Gori (ed.), *Second Workshop on Neural Networks for Speech Processing*, pp. 5–17, Trieste, Italy, 1993. Università di Firenze, Edizioni LINT Trieste S.r.l.