

# Noise Injection into Inputs in Back-Propagation Learning

Kiyotoshi Matsuoka

**Abstract**—Back-propagation can be considered a nonlinear regression technique, allowing a nonlinear neural network to acquire an input/output (I/O) association using a limited number of samples chosen from a population of input and output patterns. A crucial problem on back-propagation is its generalization capability. A network successfully trained for given samples is not guaranteed to provide desired associations for untrained inputs as well. Concerning this problem some authors showed experimentally that the generalization capability could remarkably be enhanced by training the network with noise injected inputs. This paper mathematically explains why and how the noise injection to inputs has such an effect.

## I. INTRODUCTION

**B**ACK-PROPAGATION [1] may be the most widely used method among various learning methods for neural networks. The method can be considered a kind of nonlinear regression technique, allowing a nonlinear, layered network to acquire an input/output (I/O) relation using a set of sample patterns. During the learning process, the strengths of connections between neural units are iteratively modified so that an error between the actual and desired outputs of the network will be minimized.

A crucial problem on back-propagation is its generalization capability. Usually, training patterns used for learning are only a limited number of samples picked up from a population of input and output patterns. There is no guarantee that the network successfully trained on the given samples provides desired I/O associations for untrained patterns as well. Concerning this problem some authors showed experimentally that adding some noise to input patterns in back-propagation learning remarkably enhanced the generalization capability of the resultant networks [2], [3]. This paper mathematically discusses the reason why the noise injection into input has such an effect.

## II. THE STANDARD BACK-PROPAGATION

The neural network considered in this paper is the familiar layered network with no recurrent connection. The constituent units (neurons) in it, labelled as  $1, \dots, N$ , are grouped into three types: input units  $\{1, \dots, N_I\}$ , hidden units  $\{N_I + 1, \dots, N_I + N_H\}$ , and output units  $\{N_I + N_H, \dots, N_I + N_H + N_O\}$  ( $N = N_I + N_H + N_O$ ). Input unit  $i$  receives external input  $s_i$  ( $i = 1, \dots, N_I$ ) to the network

and transmits it to other units, so output  $x_i$  of unit  $i$  is given by

$$x_i = s_i \quad (i = 1, \dots, N_I). \quad (1)$$

The output of hidden or output unit  $i$  is given by

$$x_i = g \left( \sum_j w_{ij} x_j + \theta_i \right) \quad (i = N_I + 1, \dots, N) \quad (2)$$

where  $w_{ij}$  represents the strength of the connection from unit  $j$  to unit  $i$ , and  $\theta_i$  is a bias;  $g(\cdot)$  is a nonlinear function that specifies the output of each unit in terms of its total input. The following logistic function is usually employed for it:

$$g(z) = 1 / \{1 + \exp(-z)\}. \quad (3)$$

The hidden units have no direct interaction with the external world while the output units provide the network's output  $y_i$  ( $i = 1, \dots, N_O$ ) to the external:

$$y_i = x_{N_I + N_H + i}. \quad (4)$$

We write the overall input-output mapping of the network as

$$\mathbf{y} = \mathbf{f}(\mathbf{W}, \boldsymbol{\theta}; \mathbf{s}) \quad (5)$$

or more simply  $\mathbf{y} = \mathbf{f}(\mathbf{s})$ , where  $\mathbf{s} = [s_1, \dots, s_{N_I}]^T$ ,  $\mathbf{y} = [y_1, \dots, y_{N_O}]^T$ ,  $\mathbf{W} = [w_{ij}]$ ,  $\boldsymbol{\theta} = [\theta_{N_I+1}, \dots, \theta_N]^T$ , and  $\mathbf{f} = [f_1, \dots, f_{N_O}]^T$ .

Let  $\hat{\mathbf{s}}^{(m)}$  and  $\hat{\mathbf{y}}^{(m)}$  ( $m = 1, \dots, M$ ) be a given set of pairs of input and output to be learned. In order to obtain a network that produces output  $\hat{\mathbf{y}}^{(m)}$  for input  $\hat{\mathbf{s}}^{(m)}$ , we determine the values of  $\mathbf{W}$  and  $\boldsymbol{\theta}$  such that they will minimize the following error function:

$$Q(\mathbf{W}, \boldsymbol{\theta}) = \sum_m \left| \hat{\mathbf{y}}^{(m)} - \mathbf{f}(\hat{\mathbf{s}}^{(m)}) \right|^2 \quad (6)$$

where  $|\cdot|$  denotes the Euclidean norm of a vector. A simple way for minimizing  $Q$  is to give arbitrary values to  $\mathbf{W}$  and  $\boldsymbol{\theta}$  initially and then to iteratively modify them in the opposite direction of the gradient of the error function, i.e., the modification of  $\mathbf{W}$  and  $\boldsymbol{\theta}$  should be

$$\Delta w_{ij} = -\eta \partial Q / \partial w_{ij} \quad \Delta \theta_i = -\eta \partial Q / \partial \theta_i \quad (7)$$

where  $\eta$  is a small positive constant. The calculations of  $\partial Q / \partial w_{ij}$  and  $\partial Q / \partial \theta_i$  can be performed efficiently by the well-known back-propagation algorithm [1].

If the previous iterative calculations have attained the minimum of  $Q$ , the resultant network might be optimal for the

Manuscript received June 30, 1990; revised May 22, 1991.

The author is with the Division of Control Engineering, Kyushu Institute of Technology, Tobata, Kitakyushu, 804 Japan.

IEEE Log Number 9104118.

sample patterns used in the learning. However, there is no guarantee that the network is also capable to provide reasonable associations for untrained patterns because the training patterns are usually only a small number of samples picked up from a population of an infinite number of patterns. For the network to have some generalization ability for unrepresented patterns we introduce in the next section a new performance function to be minimized.

### III. NOISE INJECTION TO INPUTS

It is an ill-posed question how to optimize  $\mathbf{W}$  and  $\theta$  for every possible pattern because we rarely have information about unrepresented patterns. So, in order to define the problem specifically we need to make some assumption on the property of the network we want to build. The principle proposed here is:

The output of the network after learning should be as insensitive as possible to input variation, as long as the error (6) is within a reasonable bound.

An essentially equivalent statement can be seen in [4]: "Similar inputs (should) tend to produce similar outputs even if they are not trained before." Obviously, these claims presuppose some smoothness of the mapping to be realized by the network.

We formulate the above principle mathematically as follows. Let  $\mathbf{s}^{(m)}$  be an input pattern similar to one training input  $\hat{\mathbf{s}}^{(m)}$ , that is

$$\mathbf{s}^{(m)} = \hat{\mathbf{s}}^{(m)} + \mathbf{d} \quad (8)$$

where  $\mathbf{d} = [d_1, \dots, d_{N_I}]^T$  and  $d_i$  are small random variables. Then, the variation of the network's output due to the disturbance,  $\mathbf{d}$ , is

$$\begin{aligned} \delta \mathbf{y}^{(m)} &= \mathbf{f}(\hat{\mathbf{s}}^{(m)} + \mathbf{d}) - \mathbf{f}(\hat{\mathbf{s}}^{(m)}) \\ &\sim \left\{ \frac{\partial \mathbf{f}(\hat{\mathbf{s}}^{(m)})}{\partial \mathbf{s}} \right\} \mathbf{d} \end{aligned} \quad (9)$$

where  $\partial \mathbf{f} / \partial \mathbf{s} = [\partial f_i / \partial s_j]$ . We here assume that  $\mathbf{d}$  is a random vector independent of  $\hat{\mathbf{s}}^{(m)}$  with the following mean and variance:

$$\langle \mathbf{d} \rangle = \mathbf{0}, \quad \langle \mathbf{d} \mathbf{d}^T \rangle = \sigma^2 \mathbf{E}_{N_I} \quad (10)$$

where  $\mathbf{E}_{N_I}$  is the identity matrix of size  $N_I$  and  $\langle \cdot \rangle$  denotes the statistical expectation of the argument.

We define the sensitivity  $R$  of the network as the mean ratio of the variances of  $|\delta \mathbf{y}^{(m)}|$  and  $|\mathbf{d}|$ :

$$\begin{aligned} R(\mathbf{w}, \theta) &= \sum_m \left\langle \frac{|\delta \mathbf{y}^{(m)}|^2}{|\mathbf{d}|^2} \right\rangle \\ &= \sum_m \left\langle \frac{\delta \mathbf{y}^{(m)T} \delta \mathbf{y}^{(m)}}{\mathbf{d}^T \mathbf{d}} \right\rangle. \end{aligned} \quad (11)$$

According to (9) and (10),  $R$  can be approximated as

$$\begin{aligned} R(\mathbf{W}, \theta) &\sim \sum_m \left\langle \frac{\mathbf{d}^T \left( \frac{\partial \mathbf{f}(\hat{\mathbf{s}}^{(m)})}{\partial \mathbf{s}} \right)^T}{\left( \frac{\partial \mathbf{f}(\hat{\mathbf{s}}^{(m)})}{\partial \mathbf{s}} \right) \mathbf{d}} \right\rangle \\ &\quad \cdot \left( \frac{\partial \mathbf{f}(\hat{\mathbf{s}}^{(m)})}{\partial \mathbf{s}} \right) \mathbf{d} \Bigg/ \langle \mathbf{d}^T \mathbf{d} \rangle \end{aligned}$$

$$\begin{aligned} &= \sum_m \left\langle \frac{\text{trace} \left\{ \left( \frac{\partial \mathbf{f}(\hat{\mathbf{s}}^{(m)})}{\partial \mathbf{s}} \right) \cdot \mathbf{d} \mathbf{d}^T \left( \frac{\partial \mathbf{f}(\hat{\mathbf{s}}^{(m)})}{\partial \mathbf{s}} \right)^T \right\}}{\text{trace} \mathbf{d} \mathbf{d}^T} \right\rangle \\ &= \sum_m \sigma^2 \text{trace} \left\{ \left( \frac{\partial \mathbf{f}(\hat{\mathbf{s}}^{(m)})}{\partial \mathbf{s}} \right) \cdot \left( \frac{\partial \mathbf{f}(\hat{\mathbf{s}}^{(m)})}{\partial \mathbf{s}} \right)^T \right\} / (N_I \sigma^2) \\ &= \sum_m \left\| \frac{\partial \mathbf{f}(\hat{\mathbf{s}}^{(m)})}{\partial \mathbf{s}} \right\|^2 / N_I \end{aligned} \quad (12)$$

where  $\|\cdot\|$  denotes the Euclidean norm of a matrix defined as  $\|\mathbf{A}\|^2 = \text{trace} \mathbf{A} \mathbf{A}^T = \sum_{i,j} a_{ij}^2$  ( $\mathbf{A} = [a_{ij}]$ ). Notice that (12) does not involve  $\sigma^2$ .

It is obvious that the smaller  $R$  is the less sensitive the network output is to the disturbance in the input pattern. So, if the previously mentioned principle is valid, it may be reasonable to make  $Q$  and  $R$  be as small as possible simultaneously. This consideration leads to the following performance function to be minimized:

$$\begin{aligned} P(\mathbf{W}, \theta) &= Q(\mathbf{W}, \theta) + aR(\mathbf{W}, \theta) \\ &\sim \sum_m \left\{ \left| \mathbf{y}^{(m)} - \mathbf{f}(\hat{\mathbf{s}}^{(m)}) \right|^2 + (a/N_I) \right. \\ &\quad \cdot \left. \left\| \frac{\partial \mathbf{f}(\hat{\mathbf{s}}^{(m)})}{\partial \mathbf{s}} \right\|^2 \right\} \quad (a > 0) \end{aligned} \quad (13)$$

where  $a$  is a weighting parameter.

In order to (locally) minimize this performance function by the steepest decent method we have to calculate  $\partial R / \partial w_{ij}$  and  $\partial R / \partial \theta_i$ . ( $\partial Q / \partial w_{ij}$  and  $\partial Q / \partial \theta_i$  can be obtained by the standard back-propagation.) Although those calculations are not so difficult, we can take an interesting alternative approach as follows. If we introduce random noise  $\mathbf{n} = [n_1, \dots, n_{N_I}]^T$  whose mean and covariance are given as

$$\langle \mathbf{n} \rangle = \mathbf{0}, \quad \langle \mathbf{n} \mathbf{n}^T \rangle = \varepsilon \mathbf{E}_{N_I} \quad (\varepsilon = a/N_I) \quad (14)$$

then (13) can be rewritten as

$$\begin{aligned} P(\mathbf{W}, \theta) &\sim \sum_m \left\{ \left| \hat{\mathbf{y}}^{(m)} - \mathbf{f}(\hat{\mathbf{s}}^{(m)}) \right|^2 + \varepsilon \left\| \frac{\partial \mathbf{f}(\hat{\mathbf{s}}^{(m)})}{\partial \mathbf{s}} \right\|^2 \right\} \\ &= \sum_m \left\{ \left| \hat{\mathbf{y}}^{(m)} - \mathbf{f}(\hat{\mathbf{s}}^{(m)}) \right|^2 \right. \\ &\quad \left. + \left\langle \left| \left( \frac{\partial \mathbf{f}(\hat{\mathbf{s}}^{(m)})}{\partial \mathbf{s}} \right) \mathbf{n} \right|^2 \right\rangle \right\} \\ &= \sum_m \left\langle \left| \hat{\mathbf{y}}^{(m)} - \mathbf{f}(\hat{\mathbf{s}}^{(m)}) - \left( \frac{\partial \mathbf{f}(\hat{\mathbf{s}}^{(m)})}{\partial \mathbf{s}} \right) \mathbf{n} \right|^2 \right\rangle \\ &\sim \left\langle \sum_m \left| \hat{\mathbf{y}}^{(m)} - \mathbf{f}(\hat{\mathbf{s}}^{(m)} + \mathbf{n}) \right|^2 \right\rangle. \end{aligned} \quad (15)$$

Equation (15) implies that minimizing  $P$  is virtually equivalent to training the network so as to minimize the error between the desired output  $\hat{y}^{(m)}$  and the actual output produced by noisy input  $\hat{s}^{(m)} + n$ . Obviously it can be achieved by the standard back-propagation without any modification. We have only to add the noise  $n$  to input.

Plaut, Nowlan, and Hinton [2] and Sietsma and Dow [3] found experimentally that adding some noise to the input applied to the back-propagation network had a remarkable effect on the generalization capability of the network. The above mathematical derivation reveals that the noise injected in the input contributes to reduce the sensitivity of the network to variation in the input.

The magnitude of noise, or  $\epsilon$ , can be determined in the following way. In most situations, an allowable bound can be designated on the error, say

$$Q(W, \theta) \leq \phi. \quad (16)$$

Since  $a$  in (13), which is a weighting parameter for the network's sensitivity, is proportional to  $\epsilon$ , we should let  $\epsilon$  be as large as possible, as long as  $Q(W, \theta)$  is smaller than  $\phi$ . In order to achieve this we monitor  $Q(W, \theta)$  during the learning process and control the magnitude of  $\epsilon$  as

- 1) when  $Q(W, \theta) > \phi$ , then decrease  $\epsilon$  by  $\Delta\epsilon$
- 2) when  $Q(W, \theta) < \phi$ , then increase  $\epsilon$  by  $\Delta\epsilon$ .

$\Delta\epsilon$  should be sufficiently small. By this rule, the value of  $Q(W, \theta)$  approaches  $\phi$  after a time of learning, while  $P(W, \theta)$  continues to decrease toward one of its local minima.

#### IV. EXAMPLES

Here, we shall show a couple of simple examples. In the first example a three-layer network with two input units, five hidden units, and one output unit was trained for three pairs of patterns:

$$\begin{aligned} \hat{s}^{(1)} &= [0, 0]^T, \hat{y}^{(1)} = 0 & \hat{s}^{(2)} &= [0.5, 0.5]^T, \hat{y}^{(2)} = 0.5 \\ \hat{s}^{(3)} &= [1, 1]^T, \hat{y}^{(3)} = 0. \end{aligned}$$

There were no direct connection from the input layer to the output layer. The standard back-propagation learning were executed with and without adding noises to input. Noises  $n_i$  ( $i = 1, 2$ ) were chosen randomly from the uniform distribution between  $-0.5$  and  $+0.5$ . Typical results are shown by the contour maps in Fig. 1. Although the learning was successful in both cases (i.e., the error converged to zero), the mapping obtained by the noiseless inputs is usually somewhat distorted as shown in Fig. 1(a), while the noisy inputs always produced apparently simpler mappings as shown in Fig. 1(b). Since every input vector,  $\hat{s}^{(m)}$ , lies on line  $s_1 = s_2$ , it might be reasonable to expect the mapping after learning to be flat in the direction perpendicular to that line. Only the result with noise injection meets this expectation.

In the second example the following patterns (exclusive OR) were applied to the same network as previously:

$$\begin{aligned} \hat{s}^{(1)} &= [0, 0]^T, \hat{s}^{(2)} = [1, 1]^T, \hat{y}^{(1)} = \hat{y}^{(2)} = 0, \\ \hat{s}^{(3)} &= [0, 0]^T, \hat{s}^{(4)} = [1, 1]^T, \hat{y}^{(3)} = \hat{y}^{(4)} = 1. \end{aligned}$$

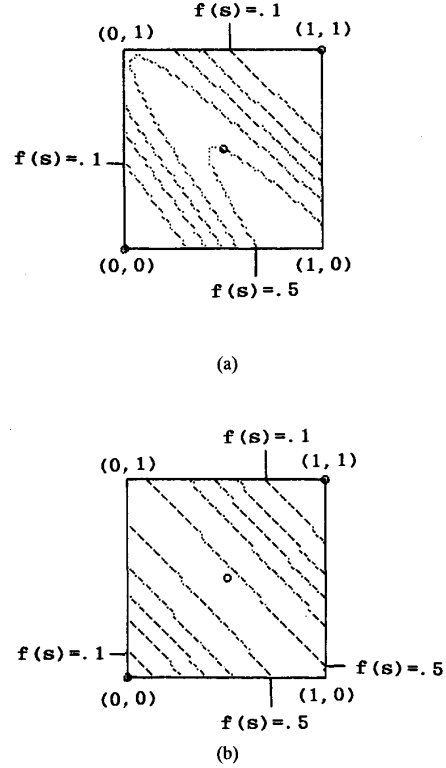


Fig. 1. Results of the learning of the first example. (a) Using noiseless input. (b) Using noisy input.

The mappings obtained by the learning with noiseless and noisy inputs are shown in Fig. 2. From the symmetric configuration of the training vectors with respect to  $s = [0.5, 0.5]^T$ , it seems reasonable to expect  $y = 0.5$  at  $s = [0.5, 0.5]^T$ . It was always achieved only with noisy input (Fig. 2(b)).

The validity of the controlling strategy for the noise magnitude,  $\epsilon$ , was tested using the same training patterns as the above two examples. In this case the noise was randomly chosen from the uniform distribution between  $-\delta$  and  $+\delta$ , in which  $\delta(\propto \epsilon)$  varied according to the controlling rule with  $\phi = 0.5$ . Mappings by the obtained networks are shown in Figs. 3(a) and (b), being similar to Fig. 1(b) and Fig. 2(b), respectively.

#### V. CONCLUSION

If the network has no hidden unit and  $g$  is a linear function (i.e.,  $g(z) = z$  and  $\theta_i = 0$  in (2)), we can show that the back-propagation with noise-injected inputs produces a network which embodies the optimal associative mapping proposed by Kohonen [5], as follows. The I/O relation of the linear network is represented by

$$y = f(s) = Ws \quad (17)$$

where  $W = [w_{ij}]$  and  $w_{ij}$  is the weight of the connection from input unit  $j$  to output unit  $i$ . Since  $\partial f(s)/\partial s = W$ , the performance function (13) reads

$$P(W) = \sum_m \left| \hat{y}^{(m)} - W\hat{s}^{(m)} \right|^2 + \epsilon \|W\|^2$$

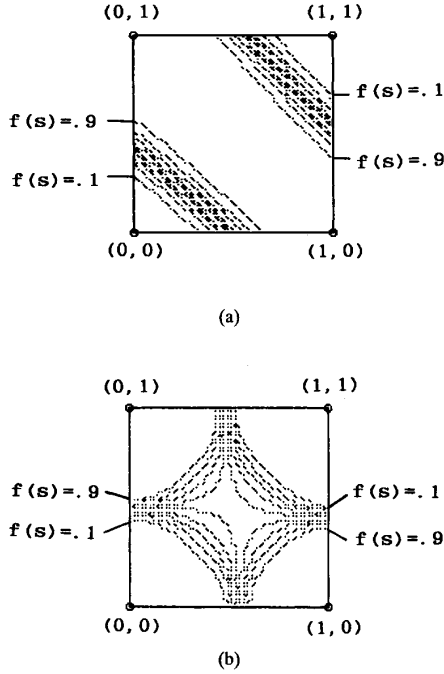


Fig. 2. Results of the learning of the second example. (a) Using noise less input. (b) Using noisy input.

$$= \|Y - WS\|^2 + \varepsilon \|W\|^2 \quad (18)$$

where  $S = [\hat{s}^{(1)}, \dots, \hat{s}^{(M)}]$  and  $Y = [\hat{y}^{(1)}, \dots, \hat{y}^{(M)}]$ . During the learning process  $W$  is iteratively modified as

$$\Delta W = -\eta \partial P(W) / \partial W = -2\eta (WSS^T - YS^T) - 2\eta \varepsilon W. \quad (19)$$

According to (19), (18) has only one local minimum, i.e., global minimum

$$W = YS^T (SS^T + \varepsilon E_{N_I})^{-1}. \quad (20)$$

Hence,  $W$  approaches this value with time for any initial values if the learning coefficient  $\eta$  is appropriately small. (20) converges with  $\varepsilon \rightarrow 0$  ( $\varepsilon \neq 0$ ) to

$$W = YS^+ \quad (21)$$

where  $(\cdot)^+$  represents the pseudoinverse of  $(\cdot)$ . This implies that if we choose  $\varepsilon$  to be very small,  $W$  becomes nearly (21), being the optimal associative mapping proposed by Kohonen [5]. The fact that the noise injection method provides a nearly optimal mapping in the linear network might encourage the present approach.

The above linear analysis suggests another interpretation of the noise injection method. The global minima of the error function  $Q(W) = \|Y - WS\|^2$  are attained by any  $W$  satisfying

$$-WSS^T + YS^T = 0$$

or

$$W = YS^+ + Z(E_{N_I} - SS^+) \quad (22)$$

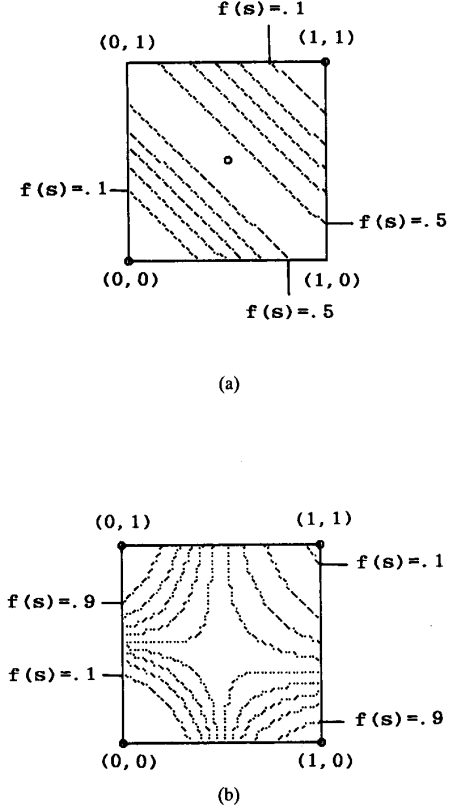


Fig. 3. Results of the learning with controlling the variance of noise.

where  $Z$  is an arbitrary  $N_O \times N_I$  matrix. If  $SS^T$  is nonsingular (in other words if  $\hat{s}^{(1)}, \dots, \hat{s}^{(M)}$  span the entire of the input space), then (22) becomes (21) because  $SS^+ = E_{N_I}$  holds. On the other hand, if  $SS^T$  is singular (in other words, if  $\hat{s}^{(1)}, \dots, \hat{s}^{(M)}$  span a strict subspace in the input space), then  $W$  cannot uniquely be determined. Even in the latter case, if we determine  $W$  such that it takes the minimum norm among  $W$  given by (22), then we have (21) uniquely. Thus, when the noise injection is applied to the linear network, it has a role to determine the optimal value of  $W$  uniquely in any case.

It is often pointed out that, in networks with hidden layers, reducing the number of hidden units to an appropriate size is effective for improving the generalization capability of the network [2], [6]. This is because the dimension of the image of the mapping is lowered by reducing the number of hidden units. In what follows, we show that the noise addition is useful even after the hidden units size is appropriately reduced.

Consider a three-layer, linear network with connection matrix  $W_1$  from the input layer to the hidden layer and connection matrix  $W_2$  from the hidden layer to the output layer. In this case, the overall I/O mapping is

$$y = W_2 W_1 s \quad (23)$$

and the error function reads

$$Q(W_1, W_2) = \sum_m |W_2 W_1 \hat{s}^{(m)} - \hat{y}^{(m)}|^2$$

$$= \|\mathbf{W}_2 \mathbf{W}_1 \mathbf{S} - \mathbf{Y}\|^2 \quad (24)$$

Let  $\mathbf{W}_1 = \mathbf{V}_1$  and  $\mathbf{W}_2 = \mathbf{V}_2$  minimize  $Q$ . Then, in general

$$\mathbf{W}_1 = \mathbf{V}_1 + \mathbf{Z}(\mathbf{E}_{N_I} - \mathbf{S}\mathbf{S}^+), \quad \mathbf{W}_2 = \mathbf{V}_2 \quad (25)$$

(where  $\mathbf{Z}$  is an arbitrary matrix with proper size) give the minimum of  $Q$  because the network generate the same output for  $\mathbf{s}^{(m)}$ :

$$\mathbf{W}_1 \mathbf{S} = \{\mathbf{V}_1 + \mathbf{Z}(\mathbf{E}_{N_I} - \mathbf{S}\mathbf{S}^+)\} \mathbf{S} = \mathbf{V}_1 \mathbf{S} \quad (26)$$

or

$$\mathbf{W}_1 \hat{\mathbf{s}}^{(m)} = \mathbf{V}_1 \hat{\mathbf{s}}^{(m)} \quad (m = 1, \dots, M). \quad (27)$$

This means that even if the dimension of  $\text{Im } \mathbf{W}_2 \mathbf{W}_1$  is lowered by reducing the number of hidden units, the mapping minimizing  $Q(\mathbf{W}_1, \mathbf{W}_2)$  or

$$\mathbf{W}_2 \mathbf{W}_1 = \mathbf{V}_2 \mathbf{V}_1 + \mathbf{V}_2 \mathbf{Z}(\mathbf{E}_{N_I} - \mathbf{S}\mathbf{S}^+) \quad (28)$$

is not necessarily unique. By applying the noise injection method to such as case, the mapping  $\mathbf{W}_2 \mathbf{W}_1$  is uniquely determined such that a weighted sum of the error function and  $\|\mathbf{W}_2 \mathbf{W}_1\|^2$  are minimized. Thus, even after reducing the hidden units to an appropriate number, noise injection to inputs can yield a remarkable effect on the generalization; one can see a good example in Seitsma and Dow [3].

When the number of training patterns is not less than the dimension of the input space, it will rarely happen that the training input vectors lie on the subspace in the input space. However, it is by no means rare that the training input vectors are located near a subspace because the patterns to be learned are more or less similar to each other in actual applications. In such a situation, the error function  $Q(\mathbf{W})$  takes the shape of a long ravine as shown in Fig. 4, and its minimum point ( $\bar{\mathbf{W}}$  in the figure) becomes quite sensitive to the samples that have happened to be used for learning. This problem is the collinearity problem which often occurs in linear regression analysis. A well-known solution to this problem is to minimize  $P(\mathbf{W})$  with finite  $\varepsilon$  even when  $Q(\mathbf{W})$  has a unique minimum point (the Ridge estimation). Then, the optimal mapping is obtained as a nearest point in the ravine from the origin ( $\tilde{\mathbf{W}}$  in Fig. 4). In Section III we implemented this well-known technique in linear regression analysis into the back-propagation learning for multilayered networks with nonlinear units.

The second term in (19),  $-2\eta\varepsilon w_{ij}$ , which cause  $w_{ij}$  to decay in proportion to current  $w_{ij}$ , can be interpreted as a kind of forgetting. In the case of the two-layer linear network, noise injection into input is equivalent to adding the decay term to the standard back-propagation. The idea of incorporating this kind of decay terms into the learning of nonlinear multilayer networks were suggested in [2], [7], [8]. As for the present method, however, the decay term for  $w_{ij}$  in a nonlinear, multilayered network takes a more complicated

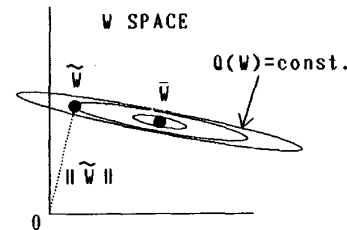


Fig. 4. Optimal weight for the error function and that for the proposed performance function.

form of  $-\eta\varepsilon \sum_m \partial \|f(\mathbf{s}^{(m)})\|^2 / \partial \mathbf{s} / \partial w_{ij}$ , depending on  $\mathbf{s}^{(m)}$  ( $m = 1, \dots, M$ ) and other connection weights.

Finally, it should be noted that the acquisition of generalization capability by noise injection into inputs relies on the assumption that the (unknown) mapping from the input space to the output space should be smooth. When this assumption is valid, the method can enhance the generalization ability of the back-propagation learning. On the other hand, when this condition is not the case, the noise injection might result in poor generalization capability.

#### REFERENCES

- [1] D.E. Rumelhart, G.E. Hinton, and R.J. Williams, "Learning internal representations by error propagation," *Parallel Distributed Processing*, D.E. Rumelhart and J.L. McClelland and the PDP Research Group, Eds. Cambridge, MA: MIT Press, pp. 318–362, 1986.
- [2] D.C. Plaut, S.J. Nowlan, and G.E. Hinton, "Experiments on learning by back-propagation," Tech. Rep. CMU-CS-86-126, 1986.
- [3] J. Sietsma and R.J.F. Dow, "Neural network pruning—Why and how," in *Proc. IEEE Int. Conf. Neural Networks*, vol. I, 1988, pp. 325–333.
- [4] E. Erusli and H. Tolle, "Learning control structures with neuron-like associative memory systems," *Organization of Neural Networks*, W. von Seelen, G. Shaw, and U.M. Leinhos, Eds. Weinheim, Germany: VCH Publishers, 1988, pp. 369–393.
- [5] T. Kohone, *Self-Organization and Associative Memory*, second ed. New York: Springer, 1987.
- [6] S.K. Kung and J.N. Hwang, "An algebraic projection analysis for optimal hidden units size and learning rates in back-propagation learning," in *Proc. IEEE Int. Conf. Neural Networks*, vol. I, 1988, pp. 363–370.
- [7] M. Ishikawa, "A structural learning algorithm with forgetting of link weights," IEICE Tech. Rep., MBE-88-144, pp. 143–148 (in Japanese), 1988.
- [8] S.J. Hanson and L.Y. Pratt, "Comparing biases for minimal network construction with back-propagation," *Advances in Neural Information Processing Systems*, D.S. Touretzky, Ed. Palo Alto: Morgan Kaufmann, 1989, pp. 177–185.



**Kiyotoshi Matsuoka** received the B.E. degree in mechanical engineering in 1971 from Kyushu Institute of Technology, Kitakyushu, and the M.S. and Ph.D. degrees in mechanical engineering in 1973 and 1976, respectively, from the University of Tokyo, Tokyo.

From 1976 to 1978 he was a Lecturer at the Department of Control Engineering, Kyushu Institute of Technology, Kitakyushu, Japan. He is currently, a Professor there. His research interests include neural networks, bioengineering, and artificial intelligence.