



ELSEVIER

Neurocomputing 000 (2001) 000–000

NEUROCOMPUTING

www.elsevier.com/locate/neucom

1

Training radial basis neural networks with the extended Kalman filter

3

Dan Simon

5

*Department of Electrical and Computer Engineering, Cleveland State University, Stilwell Hall
Room 332, 1960 East 24th Street, Cleveland, OH 44115-2425, USA*

7

Received 26 February 2000; accepted 24 May 2001

Abstract

Radial basis function (RBF) neural networks provide attractive possibilities for solving signal processing and pattern classification problems. Several algorithms have been proposed for choosing the RBF prototypes and training the network. The selection of the RBF prototypes and the network weights can be viewed as a system identification problem. As such, this paper proposes the use of the extended Kalman filter for the learning procedure. After the user chooses how many prototypes to include in the network, the Kalman filter simultaneously solves for the prototype vectors and the weight matrix. A decoupled extended Kalman filter is then proposed in order to decrease the computational effort of the training algorithm. Simulation results are presented on reformulated radial basis neural networks as applied to the Iris classification problem. It is shown that the use of the Kalman filter results in better learning than conventional RBF networks and faster learning than gradient descent. © 2001 Published by Elsevier Science Ltd.

Keywords: Radial basis function (RBF); Training; Optimization; Gradient descent; Kalman filter

1. Introduction

A radial basis function (RBF) neural network is trained to perform a mapping from an m -dimensional input space to an n -dimensional output space. RBFs can be used for discrete pattern classification, function approximation, signal processing, control, or any other application which requires a mapping from an input to an output. RBFs were first used for neural networks in [4]. Many RBF papers and

E-mail address: d.j.simon@csuohio.edu (D. Simon).

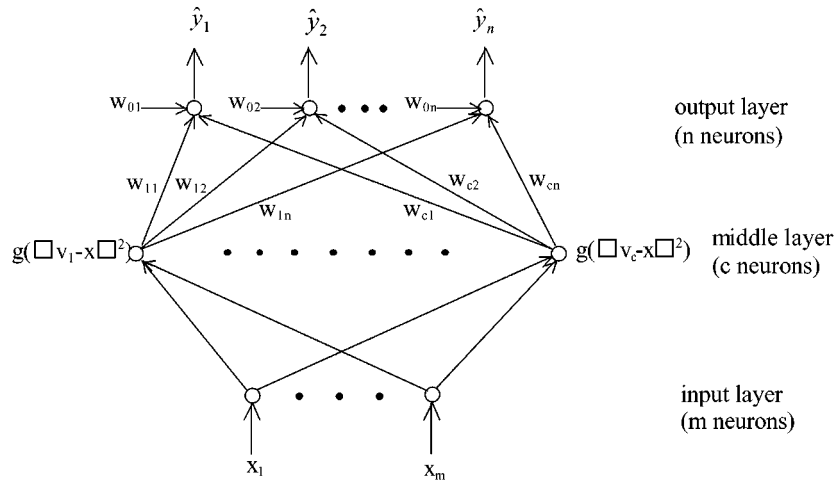


Fig. 1. Radial basis function network architecture.

1 references can be found in the recent *Neurocomputing* special issues on RBF
 2 networks [15,16].

3 An RBF consists of the m -dimensional input x being passed directly to a hidden
 4 layer. Suppose there are c neurons in the hidden layer. Each of the c neurons
 5 in the hidden layer applies an activation function which is a function of the Eu-
 6 clidean distance between the input and an m -dimensional prototype vector. Each
 7 hidden neuron contains its own prototype vector as a parameter. The output of
 8 each hidden neuron is then weighted and passed to the output layer. The outputs
 9 of the network consist of sums of the weighted hidden layer neurons. Fig. 1 shows
 a schematic of an RBF network. It can be seen that the design of an RBF requires
 several decisions, including the following:

- 11 1. How many hidden units will reside in the hidden layer (i.e., what is the value
 12 of the integer c);
- 13 2. What are the values of the prototypes (i.e., what are the values of the v vectors);
- 14 3. What function will be used at the hidden units (i.e., what is the function $g(\cdot)$);
- 15 4. What weights will be applied between the hidden layer and the output layer
 16 (i.e., what are the values of the w weights).

The performance of an RBF network depends on the number and location (in
 17 the input space) of the centers, the shape of the RBF functions at the hidden units,
 18 and the method used for determining the network weights. Some researchers have
 19 trained RBF networks by selecting the centers randomly from the training data
 20 [4]. Others have used unsupervised procedures (such as the k -means algorithm)
 21 for selecting the RBF centers [11]. Still others have used supervised procedures
 22 for selecting the RBF centers [9].

1 Several training methods separate the tasks of prototype determination and weight
2 optimization. This trend probably arose because of the quick training that could
3 result from the separation of the two tasks. In fact, one of the primary contributors
4 to the popularity of RBF networks was probably their fast training times as com-
5 pared to gradient descent training (include backpropagation). Reviewing Fig. 1,
6 it can be seen that once the prototypes are fixed and the hidden layer function $g(\cdot)$
7 is known, the network is linear in the weight parameters w . At that point training
8 the network becomes a quick and easy task that can be solved via linear least
9 squares. (This is similar to the popularity of the optimal interpolative net that is
10 due in large part to the efficient noniterative learning algorithms that are available
11 [18,19].)

12 Training methods that separate the tasks of prototype determination and weight
13 optimization often do not use the input–output data from the training set for
14 the selection of the prototypes. For instance, the random selection method and
15 the k -means algorithm result in prototypes that are completely independent of
16 the input–output data from the training set. Although this results in fast training,
17 it clearly does not take full advantage of the information contained in the training
18 set.

19 Gradient descent training of RBF networks has proven to be much more effective
20 than more conventional methods [9]. However, gradient descent training can be
21 computationally expensive. This paper extends the results of [9] and formulates a
22 training method for RBFs based on Kalman filtering. This new method proves to
23 be quicker than gradient descent training while still providing performance at the
24 same level of effectiveness.

25 Training a neural network is, in general, a challenging nonlinear optimization
26 problem. Various derivative-based methods have been used to train neural networks,
27 including gradient descent [9], Kalman filtering [20,21], and the well-known back-
28 propagation [7]. Derivative-free methods, including genetic algorithms [6], learning
29 automata [12], and simulated annealing [10], have also been used to train neural
30 networks.

31 Derivative-free methods have the advantage that they do not require the deriva-
32 tive of the objective function with respect to the neural network parameters. They
33 are more robust than derivative-based methods with respect to finding a global
34 minimum and with respect to their applicability to a wide range of objective
35 functions and neural network architectures. However, they typically tend to
36 converge more slowly than derivative-based methods. Derivative-based methods
37 have the advantage of fast convergence, but they tend to converge to local
38 minima. In addition, due to their dependence on analytical derivatives, they are
39 limited to specific objective functions and specific types of neural network
40 architectures.

41 In this paper we formulate a training method for RBFs that is based on Kalman
42 filtering. Kalman filters have been used extensively with neural networks. They
43 have been used to train multilayer perceptrons [17,20,21] and recurrent networks
44 [13,14]. They have also been used to train RBF networks, but so far their appli-
45 cation has been restricted to single-output networks with exponential functions at

1 the hidden layer [3]. In the present paper we extend the use of Kalman filters to
 2 the training of general multi-input, multi-output RBF networks.

3 For linear dynamic systems with white process and measurement noise, the
 4 Kalman filter is known to be an optimal estimator. For nonlinear systems with
 5 colored noise, the Kalman filter can be extended by linearizing the system around
 6 the current parameter estimates. This algorithm updates parameters in a way that
 7 is consistent with all previously measured data and generally converges in a few
 8 iterations. In the following sections we describe how the extended Kalman filter
 9 can be applied to RBF optimization. We demonstrate its performance on the
 10 Iris classification problem and compare it with RBF optimization using gradient
 11 descent.

12 The next section provides an overview of reformulated RBFs, and Section 3 dis-
 13 cusses how gradient-based methods can optimize the weights and prototype vectors
 14 of an RBF. Section 4 shows how an extended Kalman filter can optimize the pa-
 15 rameters of an RBF, and Section 5 proposes a modification of the Kalman filter in
 16 order to decrease the computational effort for large problems. Section 6 contains
 17 simulation results and a comparison of the Kalman filter method with the gradient
 18 descent method, and Section 7 contains some concluding remarks and suggestions
 19 for further research.

2. Reformulated radial basis functions

21 There have been a number of popular choices for the $g(\cdot)$ function at the hidden
 22 layer of RBFs [6] (see Fig. 1). The most common choice is a Gaussian function
 23 of the form

$$g(v) = \exp(-v/\beta^2), \quad (1)$$

24 where β is a real constant. Other hidden layer functions that have often been used
 25 are the thin plate spline function

$$g(v) = v \log \sqrt{v} \quad (2)$$

the multiquadric function

$$g(v) = (v^2 + \beta^2)^{1/2} \quad (3)$$

27 and the inverse multiquadric function

$$g(v) = (v^2 + \beta^2)^{-1/2}, \quad (4)$$

where β is a real constant.

29 It has recently been proposed [9] that since RBF prototypes are generally in-
 30 terpreted as the centers of receptive fields, hidden layer functions should have the
 31 following properties:

- 32 1. The response at a hidden neuron is always positive;
- 33 2. The response at a hidden neuron becomes stronger as the input approaches the
 prototype;

1 3. The response at a hidden neuron becomes more sensitive to the input as the
input approaches the prototype.

3 With these desired properties in mind, an RBF's hidden layer function should be
of the general form

$$g(v) = [g_0(v)]^{1/(1-p)} \quad (5)$$

5 and one of two sets of conditions on the so-called generator function $g_0(v)$ should
hold. If p is a real number greater than 1, then the generator function $g_0(v)$ should
satisfy the following set of conditions:

- 7
1. $g_0(v) > 0 \forall v \in (0, \infty)$;
 - 9 2. $g'_0(v) > 0 \forall v \in (0, \infty)$;
 3. $\frac{p}{p-1}[g'_0(v)]^2 - g_0(v)g''_0(v) > 0 \forall v \in (0, \infty)$.

11 The last two conditions ensure that $g'(v) < 0$ and $g''(v) > 0 \forall v \in (0, \infty)$. If p
is a real number less than 1, then the generator function $g_0(v)$ should satisfy the
following set of conditions:

- 13
1. $g_0(v) > 0 \forall v \in (0, \infty)$;
 - 15 2. $g'_0(v) < 0 \forall v \in (0, \infty)$;
 3. $\frac{p}{p-1}[g'_0(v)]^2 - g_0(v)g''_0(v) < 0 \forall v \in (0, \infty)$.

17 The last two conditions again ensure that $g'(v) < 0$ and $g''(v) > 0 \forall v \in (0, \infty)$.
It can be shown that the functions of Eqs. (1) and (4) satisfy these conditions, but
the functions of Eqs. (2) and (3) do not.

19 One new generator function that satisfies the above conditions [9] is the linear
21 function

$$g_0(v) = av + b, \quad (6)$$

23 where $a > 0$ and $b \geq 0$. If $a = 1$ and $p = 3$, then the hidden layer function reduces
to the inverse multiquadric function of Eq. (4). In this paper, we will concentrate
on hidden layer functions of the form of Eq. (5) with special emphasis on the
25 inverse multiquadric function of Eq. (4).

3. Derivative-based optimization of radial basis functions

27 The response of an RBF of the form of Fig. 1, where the hidden layer functions
 $g(\cdot)$ have the form of Eq. (5), can be written as follows:

$$\hat{y} = \begin{bmatrix} w_{10} & w_{11} & \cdots & w_{1c} \\ w_{20} & w_{21} & \cdots & w_{2c} \\ \vdots & \vdots & \vdots & \vdots \\ w_{n0} & w_{n1} & \cdots & w_{nc} \end{bmatrix} \begin{bmatrix} 1 \\ g(\|x - v_1\|^2) \\ \vdots \\ g(\|x - v_c\|^2) \end{bmatrix}. \quad (7)$$

29

- 1 We will use the following notation as shorthand for the weight matrix on the right-hand side of Eq. (7)

$$\begin{bmatrix} w_{10} & w_{11} & \cdots & w_{1c} \\ w_{20} & w_{21} & \cdots & w_{2c} \\ \vdots & \vdots & \vdots & \vdots \\ w_{n0} & w_{n1} & \cdots & w_{nc} \end{bmatrix} = \begin{bmatrix} w_1^T \\ w_2^T \\ \vdots \\ w_n^T \end{bmatrix} = W. \quad (8)$$

- 3 If we are given a training set of M desired input–output responses $\{x_i, y_i\}$ ($i = 1, \dots, M$), then we can augment M equations of the form of Eq. (7) as follows:

$$[\hat{y}_1 \quad \cdots \quad \hat{y}_M] = W \begin{bmatrix} 1 & \cdots & 1 \\ g(\|x_1 - v_1\|^2) & \cdots & g(\|x_M - v_1\|^2) \\ \vdots & \vdots & \vdots \\ g(\|x_1 - v_c\|^2) & \cdots & g(\|x_M - v_c\|^2) \end{bmatrix}. \quad (9)$$

- 5 We will introduce the following notation for the matrix on the right-hand side of Eq. (9).

$$7 \quad h_{0k} = 1 \quad (k = 1, \dots, M), \quad (10)$$

$$h_{jk} = g(\|x_k - v_j\|^2) \quad (k = 1, \dots, M), \quad (j = 1, \dots, c) \quad (11)$$

in which case we can write the matrix on the right-hand side of Eq. (9) as

$$\begin{bmatrix} h_{01} & \cdots & h_{0M} \\ h_{11} & \cdots & h_{1M} \\ \vdots & \vdots & \vdots \\ h_{c1} & \cdots & h_{cM} \end{bmatrix} = [h_1 \quad \cdots \quad h_M] = H. \quad (12)$$

- 9 In this case we can rewrite Eq. (9) as

$$\hat{Y} = WH. \quad (13)$$

- 11 Now, if we want to use gradient descent to minimize the training error, we can define the error function

$$E = \frac{1}{2} \|Y - \hat{Y}\|_F^2, \quad (14)$$

- 13 where Y is the matrix of target (desired) values for the RBF output, and $\|\cdot\|_F^2$ is the square of the Froebinius norm of a matrix, which is equal to the sum of the squares of the elements of the matrix. It has been shown [9] that in this case

$$\frac{\partial E}{\partial w_i} = \sum_{k=1}^M (\hat{y}_{ik} - y_{ik}) h_k \quad (i = 1, \dots, n), \quad (15)$$

- 15 $\frac{\partial E}{\partial v_j} = \sum_{k=1}^M 2g'(\|x_k - v_j\|^2)(x_k - v_j) \sum_{i=1}^n (y_{ik} - \hat{y}_{ik}) w_{ij} \quad (j = 1, \dots, c), \quad (16)$

- 1 where \hat{y}_{ik} is the element in the i th row and k th column of the \hat{Y} matrix of Eq. (13),
 3 RBF with respect to the rows of the weight matrix W and the prototype locations v_j
 by iteratively computing the above partials and performing the following updates:

$$\begin{aligned} w_i &= w_i - \eta \frac{\partial E}{\partial w_i} \quad (i = 1, \dots, n), \\ v_j &= v_j - \eta \frac{\partial E}{\partial v_j} \quad (j = 1, \dots, c), \end{aligned} \tag{17}$$

- 5 where η is the step size of the gradient descent method. This optimization stops
 when w_i and v_j reach local minima.

7 **4. Radial basis function optimization using the Kalman filter**

- 9 Alternatively, we can use Kalman filtering to minimize the training error. Deriva-
 11 tions of the extended Kalman filter are widely available in the literature [1,8]. In
 this section we briefly outline the algorithm and show how it can be applied to
 RBF network optimization. Consider a nonlinear finite dimensional discrete time
 system of the form

$$\begin{aligned} \theta_{k+1} &= f(\theta_k) + \omega_k, \\ y_k &= h(\theta_k) + v_k, \end{aligned} \tag{18}$$

- 13 where the vector θ_k is the state of the system at time k , ω_k is the process noise,
 15 y_k is the observation vector, v_k is the observation noise, and $f(\cdot)$ and $h(\cdot)$ are
 nonlinear vector functions of the state. Assume that the initial state θ_0 and the
 noise sequences $\{v_k\}$ and $\{\omega_k\}$ are Gaussian and independent from each other
 17 with

$$AE(\theta_0) = \bar{\theta}_0, \tag{19}$$

$$AE[(\theta_0 - \bar{\theta}_0)(\theta_0 - \bar{\theta}_0)^T] = P_0, \tag{20}$$

$$AE(\omega_k) = 0, \tag{21}$$

$$AE(\omega_k \omega_l^T) = Q \delta_{kl}, \tag{22}$$

$$AE(v_k) = 0, \tag{23}$$

$$AE(v_k v_l^T) = R \delta_{kl}, \tag{24}$$

- 19 where $AE(\cdot)$ is the expectation operator and δ_{kl} is the Kronecker delta. The problem
 addressed by the extended Kalman filter is to find an estimate $\hat{\theta}_{n+1}$ of θ_{k+1} given
 y_j ($j = 0, \dots, k$).

1 If the nonlinearities in Eq. (18) are sufficiently smooth, we can expand them
 around the state estimate $\hat{\theta}_k$ using Taylor series to obtain

$$\begin{aligned} f(\theta_k) &= f(\hat{\theta}_k) + F_k \times (\theta_k - \hat{\theta}_k) + \text{higher order terms,} \\ h(\theta_k) &= h(\hat{\theta}_k) + H_k^T \times (\theta_k - \hat{\theta}_k) + \text{higher order terms,} \end{aligned} \quad (25)$$

3 where we have introduced the notation

$$\begin{aligned} F_k &= \left. \frac{\partial f(\theta)}{\partial \theta} \right|_{\theta=\hat{\theta}_k}, \\ H_k^T &= \left. \frac{\partial h(\theta)}{\partial \theta} \right|_{\theta=\hat{\theta}_k}. \end{aligned} \quad (26)$$

5 Neglecting the higher-order terms in Eq. (25), the system in Eq. (18) can be
 approximated as

$$\begin{aligned} \theta_{k+1} &= F_k \theta_k + \omega_k + \phi_k, \\ y_k &= H_k^T \theta_k + v_k + \varphi_k, \end{aligned} \quad (27)$$

where ϕ_k and φ_k are defined as

$$\begin{aligned} \phi_k &= f(\hat{\theta}_k) - F_k \hat{\theta}_k, \\ \varphi_k &= h(\hat{\theta}_k) - H_k^T \hat{\theta}_k. \end{aligned} \quad (28)$$

7 It can be shown that the desired estimate $\hat{\theta}_n$ can be obtained by the recursion

$$\begin{aligned} \hat{\theta}_k &= f(\hat{\theta}_{k-1}) + K_k [y_k - h(\hat{\theta}_{k-1})], \\ K_k &= P_k H_k (R + H_k^T P_k H_k)^{-1}, \\ P_{k+1} &= F_k (P_k - K_k H_k^T P_k) F_k^T + Q. \end{aligned} \quad (29)$$

9 K_k is known as the Kalman gain. In the case of a linear system, it can be
 shown that P_k is the covariance matrix of the state estimation error, and the
 state estimate $\hat{\theta}_{k+1}$ is optimal in the sense that it approaches the conditional mean
 11 $AE[\theta_{k+1} | (y_0, y_1, \dots, y_k)]$ for large k . For nonlinear systems the filter is not optimal
 and the estimates are only approximately conditional means.

13 Inspired by the successful use of the Kalman filter for training non-RBF neural
 networks [14], we can apply a similar technique to the training of RBF networks. In
 15 general, we can view the optimization of the weight matrix W and the prototypes
 v_j as a weighted least-squares minimization problem, where the error vector is
 17 the difference between the RBF outputs and the target values for those outputs.
 Consider the RBF network of Fig. 1 with m inputs, c prototypes, and n outputs.
 19 We use y to denote the target vector for the RBF outputs, and $h(\hat{\theta}_k)$ to denote

1 the actual outputs at the k th iteration of the optimization algorithm.

$$\begin{aligned}
 y &= [y_{11} \quad \cdots \quad y_{1M} \quad \cdots \quad y_{n1} \quad \cdots \quad y_{nM}]^T, \\
 h(\hat{\theta}_k) &= [\hat{y}_{11} \quad \cdots \quad \hat{y}_{1M} \quad \cdots \quad \hat{y}_{n1} \quad \cdots \quad \hat{y}_{nM}]_k^T.
 \end{aligned} \tag{30}$$

3 Note that the y and \hat{y} vectors each consist of nM elements, where n is the di-
 4 mension of the RBF output and M is the number of training samples. In order to
 5 cast the optimization problem in a form suitable for Kalman filtering, we let the
 6 elements of the weight matrix W and the elements of the prototypes v_j constitute
 7 the state of a nonlinear system, and we let the output of the RBF network consti-
 8 tute the output of the nonlinear system to which the Kalman filter is applied. The
 9 state of the nonlinear system can then be represented as

$$\theta = [w_1^T \quad \cdots \quad w_n^T \quad v_1^T \quad \cdots \quad v_c^T]^T. \tag{31}$$

9 The vector θ thus consists of all $(n(c+1) + mc)$ of the RBF parameters arranged
 10 in a linear array. The nonlinear system model to which the Kalman filter can be
 11 applied is

$$\begin{aligned}
 \theta_{k+1} &= \theta_k, \\
 y_k &= h(\theta_k),
 \end{aligned} \tag{32}$$

12 where $h(\theta_k)$ is the RBF network's nonlinear mapping between its parameters and
 13 its output. In order to execute a stable Kalman filter algorithm, we need to add
 14 some artificial process noise and measurement noise to the system model [14]. So
 15 we rewrite Eq. (32) as

$$\begin{aligned}
 \theta_{k+1} &= \theta_k + \omega_k, \\
 y_k &= h(\theta_k) + v_k,
 \end{aligned} \tag{33}$$

16 where ω_k and v_k are artificially added noise processes. Now we can apply the
 17 Kalman recursion of Eq. (29). $f(\cdot)$ is the identity mapping and y_k is the target
 18 output of the RBF network. (Note that although y_k is written as a function of the
 19 Kalman iteration number k , it is actually a constant.) $h(\hat{\theta}_k)$ is the actual output
 20 of the RBF network given the RBF parameters at the k th iteration of the Kalman
 21 recursion. H_k is the partial derivative of the RBF output with respect to the RBF
 22 network parameters at the k th iteration of the Kalman recursion. F_k is the identity
 23 matrix (again, a constant even though it is written as a function of k). The Q
 24 and R matrices are tuning parameters which can be considered as the covariance
 25 matrices of the artificial noise processes ω_k and v_k , respectively. It is shown in
 26 Appendix A that the partial derivative of the RBF output with respect to the RBF
 27 network parameters is given by

$$H_k = \begin{bmatrix} H_w \\ H_v \end{bmatrix}, \tag{34}$$

1 where H_w and H_v are given by

$$H_w = \begin{bmatrix} H & 0 & \cdots & 0 \\ 0 & H & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & H \end{bmatrix}, \quad (35)$$

$$H_v = \begin{bmatrix} -w_{11}g'_{11}2(x_1 - v_1) & \cdots & -w_{11}g'_{m1}2(x_m - v_1) & \cdots \\ \vdots & \vdots & \vdots & \vdots \\ -w_{1c}g'_{1c}2(x_1 - v_c) & \cdots & -w_{1c}g'_{mc}2(x_m - v_c) & \cdots \\ \vdots & \vdots & \vdots & \vdots \\ -w_{n1}g'_{11}2(x_1 - v_1) & \cdots & -w_{n1}g'_{m1}2(x_m - v_1) \\ \vdots & \vdots & \vdots & \vdots \\ -w_{nc}g'_{1c}2(x_1 - v_c) & \cdots & -w_{nc}g'_{mc}2(x_m - v_c) \end{bmatrix}, \quad (36)$$

where H (with no subscript) is the $(c + 1) \times M$ matrix given in Eq. (12), w_{ij} is the element in the i th row and j th column of the W weight matrix (see Eq. (8)), $g'_{ij} = g'(|x_i - v_j|^2)$ (where $g(\cdot)$ is the activation function at the hidden layer), x_i is the i th input vector, and v_j is the j th prototype vector. H_w in Eq. (35) is an $n(c + 1) \times nM$ matrix, H_v in Eq. (36) is an $mc \times nM$ matrix, and H_k in Eq. (34) is an $[n(c + 1) + mc] \times nM$ matrix. Appendix B gives a derivation of $g'_{ij}(\cdot)$ for the case when $g(\cdot)$ is generated by the linear generator function. Now that we have the H_k matrix, we can execute the recursion of Eq. (29), thus using the extended Kalman filter in order to determine the weight matrix W and the prototypes v_j .

5. Decoupling the Kalman filter

13 The Kalman filter described in the previous section can be decoupled in order to
 14 save computational effort. This is similar to the decoupling that has been performed
 15 for Kalman filter training of recurrent neural networks [14]. For a large RBF
 16 network, the computational expense of the Kalman filter could be burdensome. In
 17 fact, the computational expense of the Kalman filter is on the order of AB^2 , where
 18 A is the dimension of the output of the dynamic system and B is the number of
 19 parameters. In our case, there are nM outputs and $n(c + 1) + mc$ parameters, where
 20 n is the dimension of the RBF output, M is the number of training samples, c is
 21 the number of prototypes, and m is the dimension of the RBF input. Therefore, the
 22 computational expense of the Kalman filter is on the order of $nM[n(c + 1) + mc]^2$.

23 The Kalman filter parameter vector can be decoupled by assuming that certain
 24 parameter groups interact with each other only at a second-order level. For in-
 25 stance, it can be seen from Appendix A that the H_k matrix contains a lot of zeros,
 26 showing that the interaction between various parameter groups can be neglected.
 27 In particular, the H_k matrix consists of $n + 1$ decoupled blocks. These $n + 1$ blocks
 correspond to the n sets of w weights that affect the n output components, and

1 the set of prototypes. This is intuitive because, for example, the $c + 1$ weights that
 3 $c + 1$ weights that impinge on the second component of the output (see Fig. 1). Let
 5 us use the notation that θ_k^i refers to the i th group (out of $n + 1$ total) of parameters
 estimated at time step k . Then we have

$$\begin{aligned} \theta_k^1 &= w_1 \\ &\vdots \\ \theta_k^n &= w_n \\ \theta_k^{n+1} &= [v_1^T \quad \cdots \quad v_c^T]^T. \end{aligned} \quad (37)$$

7 We will use the notation that H_k^i refers to the submatrix of H_k corresponding to
 the i th group of parameters.

$$\begin{aligned} H_k^1 &= H \\ &\vdots \\ H_k^n &= H \\ H_k^{n+1} &= H_v, \end{aligned} \quad (38)$$

9 where H (with no subscript) is the $(c + 1) \times M$ matrix given in Eq. (12), and
 11 H_v is given in Eq. (36). We will use the notation that y_k^i refers to the elements
 of the target output of the RBF network that are affected by the i th group of
 parameters.

$$\begin{aligned} y_k^1 &= [y_{11} \quad \cdots \quad y_{1M}]^T \\ &\vdots \\ y_k^n &= [y_{n1} \quad \cdots \quad y_{nM}]^T \\ y_k^{n+1} &= [y_{11} \quad \cdots \quad y_{1M} \quad \cdots \quad y_{n1} \quad \cdots \quad y_{nM}]^T. \end{aligned} \quad (39)$$

13 Similarly, we use the notation that $h^i(\hat{\theta}_{k-1})$ refers to the elements of the actual
 output of the RBF network that are affected by the i th group of parameters.

$$\begin{aligned} h^1(\hat{\theta}_{k-1}) &= [\hat{y}_{11} \quad \cdots \quad \hat{y}_{1M}]_k^T \\ &\vdots \\ h^n(\hat{\theta}_{k-1}) &= [\hat{y}_{n1} \quad \cdots \quad \hat{y}_{nM}]_k^T \\ h^{n+1}(\hat{\theta}_{k-1}) &= [\hat{y}_{11} \quad \cdots \quad \hat{y}_{1M} \quad \cdots \quad \hat{y}_{n1} \quad \cdots \quad \hat{y}_{nM}]_k^T. \end{aligned} \quad (40)$$

- 1 The decoupled Kalman recursion for the i th parameter group, modified from Eq. (29), is then given by

$$\begin{aligned}\hat{\theta}_k^i &= f(\hat{\theta}_{k-1}^i) + K_k^i [y_k^i - h^i(\hat{\theta}_{k-1}^i)], \\ K_k^i &= P_k^i H_k^i (R^i + (H_k^i)^T P_k^i H_k^i)^{-1}, \\ P_{k+1}^i &= F_k (P_k^i - K_k^i (H_k^i)^T P_k^i) F_k^T + Q^i.\end{aligned}\quad (41)$$

- 3 As before, $f(\cdot)$ is the identity mapping and F_k is the identity matrix. The above recursion executes $n+1$ times. The first n times, the recursion consists of M outputs
5 and $(c+1)$ parameters. The last time, the recursion consists of nM outputs and mc parameters. So the computational expense of the Kalman filter has been reduced
7 to the order of $nM[(c+1)^2 + (mc)^2]$. The ratio of the computational expense of the standard Kalman filter to the decoupled Kalman filter can be computed as

$$\frac{\text{Standard KF Expense}}{\text{Decoupled KF Expense}} = \frac{n^2(c+1)^2 + m^2c^2 + n(c+1)mc}{(c+1)^2 + m^2c^2}. \quad (42)$$

- 9 The computational savings will be most significant for large problems, i.e., problems where n (the dimension of the output) is large, m (the dimension of the
11 input) is large, or c (the number of prototypes) is large. Note that this complexity analysis applies only to the Kalman recursion and does not include the computational expense required for the calculation of the partial derivatives (see Appendix
13 A). Also note that for both the standard and decoupled Kalman filters, the computational expense increases linearly with M (the number of training samples).
15

6. Simulation results

- 17 In this section we describe and illustrate the use of Kalman filter training for the parameters of an RBF network. We tested the algorithms of the previous
19 sections on the classical Iris classification problem [2]. Each Iris exemplar has four features and is classified into one of three categories. The Iris data contains 50 exemplars from each category for a total of 150 patterns. We randomly divided
21 the patterns into training and test sets, each containing 25 exemplars from each category. The input data were normalized by replacing each feature value x by
23 $\bar{x} = (x - \mu_x) / \sigma_x$, where μ_x and σ_x denote the sample mean and standard deviation of this feature over the entire data set. The networks were trained to respond with
25 the target value $y_{ik} = 1$, and $y_{jk} = 0 \quad \forall j \neq i$, when presented with an input vector x_k from the i th category. The reformulated RBF networks were trained using the
27 hidden layer function of Eq. (5) with the linear generator function of Eq. (6) with $a=1$ and $b=1$. The exponential parameter p in Eq. (5) was varied between 2
29 and 4 in the simulation results presented in this section. The training algorithms were initialized with prototype vectors randomly selected from the input data, and
31 with the weight matrix W set to 0.

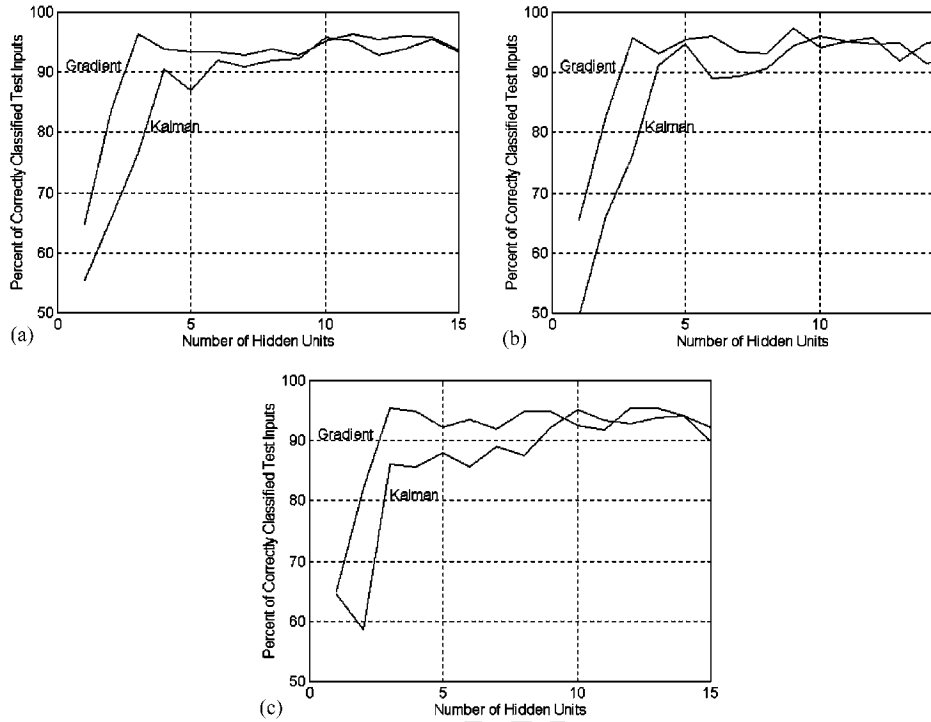


Fig. 2. Average RBF performance on the Iris test data with linear generator functions $g_0(v) = v + 1$ and $g(v) = [g_0(v)]^{1/(1-p)}$: (a) $p = 2$; (b) $p = 3$; (c) $p = 4$.

- 1 After some experimentation, it was concluded that gradient descent worked best
- 2 with $\eta = 0.01$ (see Eq. (17)). The gradient descent optimization algorithm was
- 3 terminated when the error function of Eq. (14) decreased by less than 0.1%.
- 4 The Kalman filter parameters of Eq. (29) were initialized with $P_0 = 40I$, $Q = 40I$,
- 5 and $R = 40I$, where I is the identity matrix of appropriate dimensions. The Kalman
- 6 filter recursion was terminated when the error function of Eq. (14) decreased by
- 7 less than 0.1%.
- 8 The decoupled Kalman filter parameters of Eq. (41) were initialized with $P_0^i = 40I$,
- 9 $Q^i = 40I$, and $R^i = 40I$, where $i = 1, \dots, (n + 1)$, and where I is the identity matrix
- 10 of appropriate dimensions.
- 11 The performance of each of the training methods was explored by averaging
- 12 its performance over five trials, where each trial consisted of a random selection
- 13 of training and test data. Fig. 2 depicts the performance of the RBF network on
- 14 the test data when the network was trained with gradient descent and when it was
- 15 trained with the Kalman filter. The number of hidden units in the RBF network
- 16 was varied between 1 and 15. It can be seen from the figure that, in general,
- 17 gradient descent training resulted in a better performing network than Kalman
- filter training. But as the number of hidden units increases, the performances of

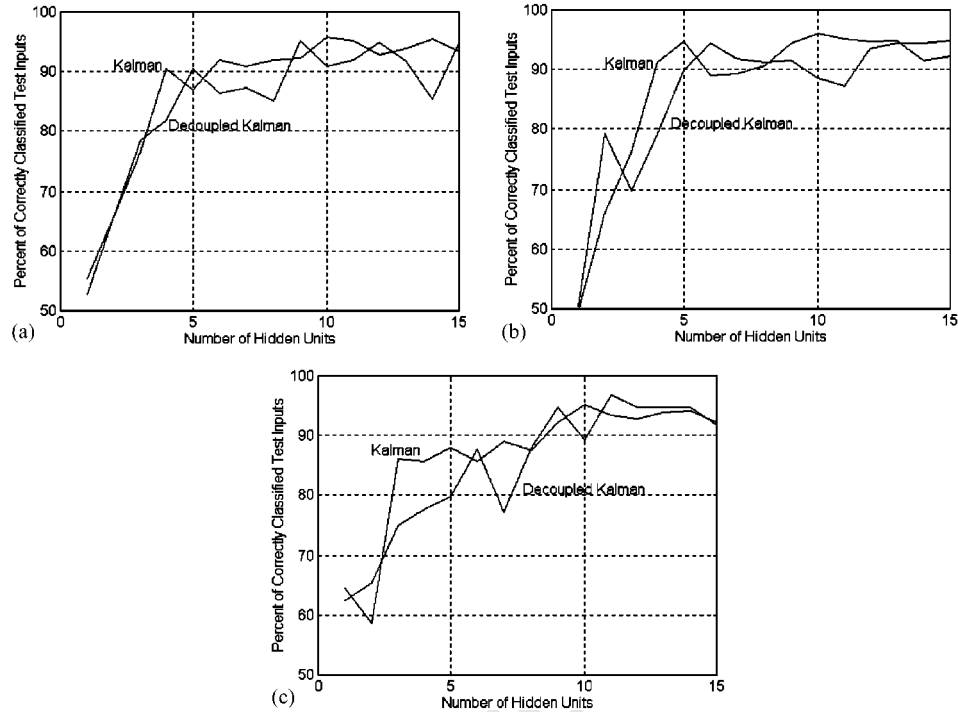


Fig. 3. Average RBF performance on the Iris test data with linear generator functions $g_0(v) = v + 1$ and $g(v) = [g_0(v)]^{1/(1-p)}$: (a) $p = 2$; (b) $p = 3$; (c) $p = 4$.

- 1 the two training algorithms are very similar. The RBF network reaches a peak
- 2 performance of about 95%. This is slightly worse than the results presented in [9].
- 3 The discrepancies could be due to differences in the training termination criteria,
- 4 or differences in other implementation details.
- 5 Fig. 3 shows the performance of the RBF network when it was trained with
- 6 the Kalman filter and when it was trained with the decoupled Kalman filter. The
- 7 number of hidden units in the RBF network was again varied between 1 and 15.
- 8 It can be seen from the figure that Kalman filter training resulted in marginally
- 9 better performance than decoupled Kalman filter training. The decoupling strategy
- 10 discussed in Section 5 degraded the performance of the RBF network only slightly.
- 11 Fig. 4 shows the number of iterations required for convergence for gradient
- 12 descent training, Kalman filter training, and decoupled Kalman filter training. It can
- 13 be seen that gradient descent training requires more iterations for convergence than
- 14 Kalman filter training (both standard and decoupled) by a full order of magnitude.
- 15 This indicates the computational superiority of Kalman filters over gradient descent.
- 16 Note that the convergence criterion is identical for all training methods—that is,
- 17 training was terminated when the error function of Eq. (14) decreased by less than
- 0.1%.

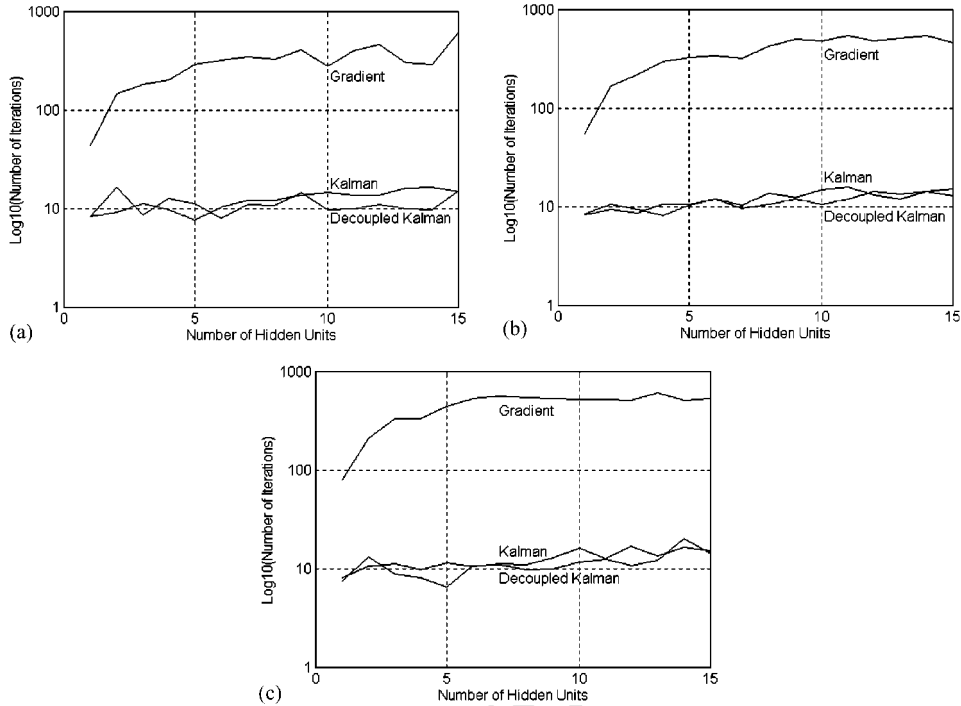


Fig. 4. Average number of iterations required for learning convergence with linear generator functions $g_0(v) = v + 1$ and $g(v) = [g_0(v)]^{1/(1-p)}$: (a) $p = 2$; (b) $p = 3$; (c) $p = 4$.

1 Fig. 5 compares the CPU time required for convergence for the three training
 2 methods. (The CPU time is measured in seconds on a Pentium III 550 MHz CPU
 3 running MATLAB.) With just one or two hidden units, the CPU time is comparable
 4 for each of the three methods. But as the number of hidden units increases above
 5 one or two, the CPU time required by gradient descent reaches a full order of
 6 magnitude greater than that required by the Kalman filters. (Recall from Figs.
 7 2 and 3 that the number of hidden units needs to be more than one or two in
 8 order to achieve good test performance.) It may be somewhat surprising from
 9 Fig. 5 that the decoupled Kalman filter requires about the same amount of CPU
 10 time as the standard Kalman filter. However, as the analysis in Section 5 shows,
 11 the computational savings achieved by the decoupled Kalman filter will be most
 12 significant for problems where the dimension of the output is large, the dimension
 13 of the input is large, or the number of prototypes is large. In our Iris classification
 14 problem, the dimension of the input is 4, the dimension of the output is 3, and the
 15 number of prototypes varies between 1 and 15. So we simply have too small of a
 16 problem to see much computational savings. Also, recall that the analysis in Section
 17 5 applies only to the Kalman recursion and does not include the computational

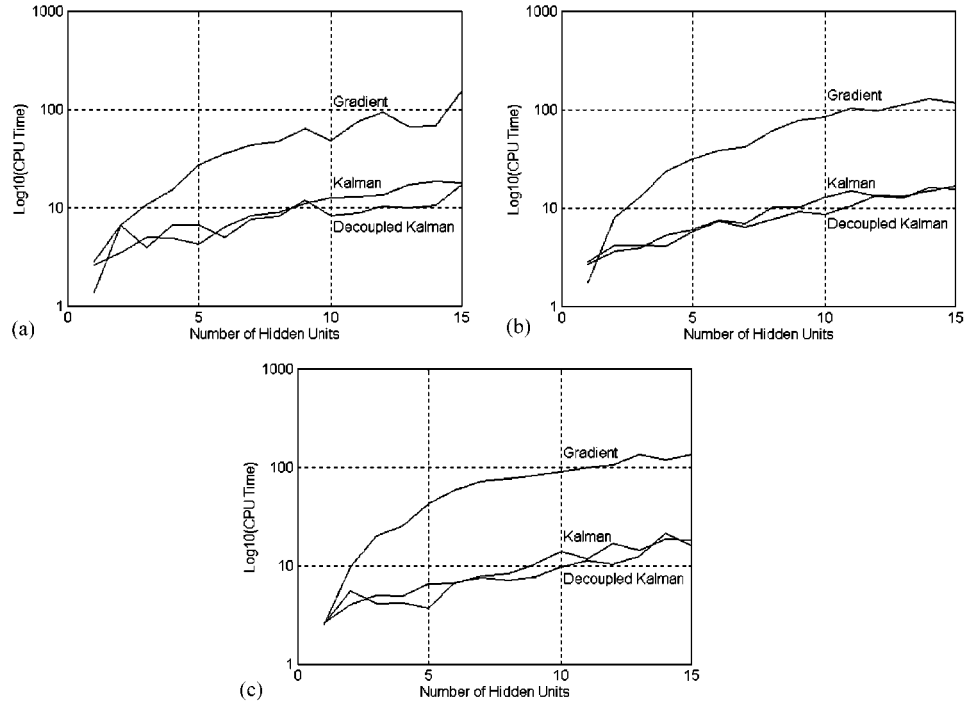


Fig. 5. Average CPU time required for learning convergence with linear generator functions $g_0(v) = v + 1$ and $g(v) = [g_0(v)]^{1/(1-p)}$: (a) $p = 2$; (b) $p = 3$; (c) $p = 4$.

1 expense required for the calculation of the partial derivatives. This further dilutes
the computational savings of the decoupled Kalman filter.

3 7. Conclusion

5 The success of a neural network architecture depends heavily on the availabil-
ity of effective learning algorithms. The theoretical strength of the Kalman filter
7 has led to its use in hundreds of technologies, and this paper demonstrates that
RBF network training is yet another fruitful application of Kalman filtering. The
9 experiments reported in this paper verify that Kalman filter training provides about
the same performance as gradient descent training, but with only a fraction of the
11 computational effort. In addition, it has been shown that the decoupled Kalman
filter provides performance on par with the standard Kalman filter while further
decreasing the computational effort for large problems.

13 Further research could focus on the application of Kalman filter training to RBF
networks with alternative forms of the generator function. In addition, the con-
15 vergence of the Kalman filter could be further improved by more intelligently
initializing the training process. (Recall that in this paper the prototypes were all

1 initialized to random input vectors and the weight matrix was initialized to zero.)
 2 Other work could focus on applying these techniques to large problems to ob-
 3 tain experimental verification of the computational savings of decoupled Kalman
 4 filter training. Additional efforts could be directed towards effective determination
 5 of the Kalman filter tuning parameters (P_0 , Q , and R). There are several objec-
 6 tives which could be addressed by further research along these lines. This paper
 7 represents just the initial steps in a direction that appears to be
 8 promising.
 9 The MATLAB m-files that were used to generate the results presented in this
 10 paper can be downloaded from the world-wide web at <http://academic.csuohio.edu/simond/rbfkalman/>. The Iris data files and the m-files for gradient descent
 11 training, Kalman filter training, and decoupled Kalman filter training are available
 12 for download and experimentation.
 13

Appendix A.

15 Use the notation \hat{Y}_j to denote the output of the RBF network given the j th
 16 training input and the current network parameters. Use the notation \hat{Y} to denote
 17 the concatenation of all M training outputs.

$$\hat{Y} = [\hat{Y}_1^T \quad \cdots \quad \hat{Y}_M^T]^T. \tag{A.1}$$

From Section 3 we can see that

$$\hat{Y} = \begin{bmatrix} \sum_{i=0}^c w_{1i} h_{i1} \\ \vdots \\ \sum_{i=0}^c w_{1i} h_{iM} \\ \vdots \\ \sum_{i=0}^c w_{ni} h_{i1} \\ \vdots \\ \sum_{i=0}^c w_{ni} h_{iM} \end{bmatrix}. \tag{A.2}$$

19 So if we define w as the first $n(c + 1)$ elements of the θ vector in Eq. (31)

$$w = [w_1^T \quad \cdots \quad w_n^T]^T \tag{A.3}$$

- 1 then the partial derivative of \hat{Y} with respect to the weights w_{ij} in the Kalman filter parameter vector of Eq. (31) can be derived as

$$\frac{\partial \hat{Y}}{\partial w} = \begin{bmatrix} \partial w_{10} \\ \vdots \\ \partial w_{1c} \\ \vdots \\ \partial w_{n0} \\ \vdots \\ \partial w_{nc} \end{bmatrix} \begin{bmatrix} \sum_{i=0}^c w_{1i} h_{i1} & \cdots & \sum_{i=0}^c w_{1i} h_{iM} & \cdots \\ & & \sum_{i=0}^c w_{ni} h_{i1} & \cdots & \sum_{i=0}^c w_{ni} h_{iM} \end{bmatrix} \quad (\text{A.4})$$

$$= \begin{bmatrix} H & 0 & \cdots & 0 \\ 0 & H & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \cdots & 0 & H \end{bmatrix}, \quad (\text{A.5})$$

- 3 where H is given in Eq. (12). From Section 3 we can also see that

$$\hat{Y} = \begin{bmatrix} w_{10} + \sum_{j=1}^c w_{1j} g(|x_1 - v_j|^2) & \cdots & w_{10} + \sum_{j=1}^c w_{1j} g(|x_M - v_j|^2) & \cdots \\ w_{n0} + \sum_{j=1}^c w_{nj} g(|x_1 - v_j|^2) & \cdots & w_{n0} + \sum_{j=1}^c w_{nj} g(|x_M - v_j|^2) \end{bmatrix}. \quad (\text{A.6})$$

So if we define v as the last mc elements of the θ vector in Eq. (31)

$$v = [v_1^T \quad \cdots \quad v_c^T]^T \quad (\text{A.7})$$

- 5 then the partial derivative of \hat{Y} with respect to the prototypes v_j in the Kalman filter parameter vector of Eq. (31) can be derived as

$$\frac{\partial \hat{Y}}{\partial v} = \begin{bmatrix} \partial v_1 \\ \vdots \\ \partial v_c \end{bmatrix} \begin{bmatrix} \sum_{j=1}^c w_{1j} g_{1j} & \cdots & \sum_{j=1}^c w_{1j} g_{Mj} & \cdots \\ & & \sum_{j=1}^c w_{nj} g_{1j} & \cdots & \sum_{j=1}^c w_{nj} g_{Mj} \end{bmatrix} \quad (\text{A.8})$$

$$= \begin{bmatrix} -w_{11}g'_{11}2(x_1 - v_1) & \cdots & -w_{11}g'_{M1}2(x_M - v_1) & \cdots \\ \vdots & \vdots & \vdots & \vdots \\ -w_{1c}g'_{1c}2(x_1 - v_c) & \cdots & -w_{1c}g'_{Mc}2(x_M - v_c) & \cdots \\ \\ -w_{n1}g'_{11}2(x_1 - v_1) & \cdots & -w_{n1}g'_{M1}2(x_M - v_1) \\ \vdots & \vdots & \vdots \\ -w_{nc}g'_{1c}2(x_1 - v_c) & \cdots & -w_{nc}g'_{Mc}2(x_M - v_c) \end{bmatrix}, \quad (\text{A.9})$$

- 1 where $g_{ij} = g(\|x_i - v_j\|^2)$, $g'_{ij} = g'(\|x_i - v_j\|^2)$ (where $g(\cdot)$ is the generator function
 2 at the hidden layer), x_i is the i th input vector, and v_j is the j th prototype vector.
 3 Now, combining Eqs. (A.5) and (A.9), we obtain Eqs. (34)–(36) for the partial
 4 derivative of the RBF network output with respect to the Kalman filter parameter
 5 vector.

Appendix B.

- 7 Consider the hidden layer function

$$g(v) = [g_0(v)]^{1/(1-p)} \quad (\text{B.1})$$

with the linear generator function

$$g_0(v) = v + \gamma^2, \quad (\text{B.2})$$

- 9 where γ is some constant. Note that $g'_0(v) = 1$. Therefore

$$g'(v) = \frac{1}{1-p} g_0(v)^{p/(1-p)} g'_0(v) \quad (\text{B.3})$$

$$= \frac{1}{1-p} (v + \gamma^2)^{p/(1-p)} \quad (\text{B.4})$$

$$= \frac{1}{1-p} g^p(v). \quad (\text{B.5})$$

- 11 Recall from Eq. (11) that $h_{jk} = g(\|x_k - v_j\|^2)$, ($k = 1, \dots, M$), ($j = 1, \dots, c$). Com-
 bining these equations, we can write

$$g'(\|x_k - v_j\|^2) = \frac{1}{1-p} g^p(\|x_k - v_j\|^2) \quad (\text{B.6})$$

$$= \frac{1}{1-p} [g_0(\|x_k - v_j\|^2)]^{p/(1-p)} \quad (\text{B.7})$$

$$= \frac{1}{1-p} h_{jk}^p. \quad (\text{B.8})$$

- 1 This provides the required expression for the derivative of the hidden layer function
that is needed to calculate the H_v matrix of Eq. (36) and thus perform the Kalman
3 recursion of Eq. (29).

8. Uncited reference

- 5 [5]

References

- 7 [1] B. Anderson, J. Moore, Optimal Filtering, Prentice-Hall, Englewood Cliffs, NJ, 1979.
8 [2] J. Bezdek, J. Keller, R. Krishnapuram, L. Kuncheva, H. Pal, Will the real Iris data please stand
9 up? IEEE Trans. Fuzzy Systems 7 (1999) 368–369.
10 [3] M. Birgmeier, A fully Kalman-trained radial basis function network for nonlinear speech
11 modeling, IEEE International Conference on Neural Networks, Perth, Western Australia, 1995,
pp. 259–264.
12 [4] D. Broomhead, D. Lowe, Multivariable functional interpolation and adaptive networks, Complex
13 Systems 2 (1988) 321–355.
14 [5] S. Chen, C. Cowan, P. Grant, Orthogonal least squares learning algorithm for radial basis function
15 networks, IEEE Trans. Neural Networks 2 (1991) 302–309.
16 [6] S. Chen, Y. Wu, B. Luk, Combined genetic algorithm optimization and regularized orthogonal
17 least squares learning for radial basis function networks, IEEE Trans. Neural Networks 10 (1999)
18 1239–1243.
19 [7] R. Duro, J. Reyes, Discrete-time backpropagation for training synaptic delay-based artificial neural
20 networks, IEEE Trans. Neural Networks 10 (1999) 779–789.
21 [8] A. Gelb, Applied Optimal Estimation, MIT Press, Cambridge, MA, 1974.
22 [9] N. Karayiannis, Reformulated radial basis neural networks trained by gradient descent, IEEE
23 Trans. Neural Networks 3 (1999) 657–671.
24 [10] S. Kirkpatrick, C.I. Gelatt, M. Vecchi, Optimization by simulated annealing, Science 220 (1983)
25 671–680.
26 [11] J. Moody, C. Darken, Fast learning in networks of locally-tuned processing units, Neural Comput.
27 1 (1989) 289–303.
28 [12] K. Narendra, M. Thathachar, Learning Automata—An Introduction, Prentice-Hall, Englewood
29 Cliffs, NJ, 1989.
30 [13] D. Obradovic, On-line training of recurrent neural networks with continuous topology adaptation,
31 IEEE Trans. Neural Networks 7 (1996) 222–228.
32 [14] G. Puskorius, L. Feldkamp, Neurocontrol of nonlinear dynamical systems with Kalman filter
33 trained recurrent networks, IEEE Trans. Neural Networks 5 (1994) 279–297.
34 [15] V.D. Sánchez, A. (Ed.), Special Issue on RBF Networks, Part I, Neurocomputing 19 (1998).
35 [16] V.D. Sánchez, A. (Ed.), Special Issue on RBF Networks, Part II, Neurocomputing 20 (1998).
36 [17] S. Shah, F. Palmieri, M. Datum, Optimal filtering algorithms for fast learning in feedforward
37 neural networks, Neural Networks 5 (1992) 779–787.
38 [18] D. Simon, Distributed fault tolerance in optimal interpolative nets, submitted for publication.
39 [19] S. Sin, R. DeFigueiredo, Efficient learning procedures for optimal interpolative nets, Neural
40 Networks 6 (1993) 99–113.
41 [20] J. Sum, C. Leung, G. Young, W. Kan, On the Kalman filtering method in neural network training
42 and pruning, IEEE Trans. Neural Networks 10 (1999) 161–166.
43 [21] Y. Zhang, X. Li, A fast U-D factorization-based learning algorithm with applications to
44 nonlinear system modeling and identification, IEEE Trans. Neural Networks 10 (1999)
45 930–938.

1
3
5
7



Dan Simon received his BS, MS, and Ph.D. degrees from Arizona State University, the University of Washington, and Syracuse University, all in Electrical Engineering. He has 14 years of industrial experience in the aerospace, automotive, biomedical, process control, and software engineering fields. He has been an Assistant Professor at Cleveland State University since 1999, where he teaches control theory, computational intelligence, and embedded systems. He has over 30 publications in referred journals and conference proceedings.

UNCORRECTED PROOF