

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/264579469>

# Recurrent Least Squares Support Vector Machines

Article in IEEE Transactions on Circuits and Systems I Fundamental Theory and Applications · July 2000

DOI: 10.1109/81.855471

CITATIONS

524

READS

327

2 authors:



Johan A.K. Suykens

[www.esat.kuleuven.be/stadius](http://www.esat.kuleuven.be/stadius)

741 PUBLICATIONS 28,059 CITATIONS

[SEE PROFILE](#)



Joos Vandewalle

KU Leuven

724 PUBLICATIONS 29,162 CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



Crop area and crop yields assessment from time series of remotely sensed images [View project](#)



Concept Inventories for Circuits and Systems courses [View project](#)

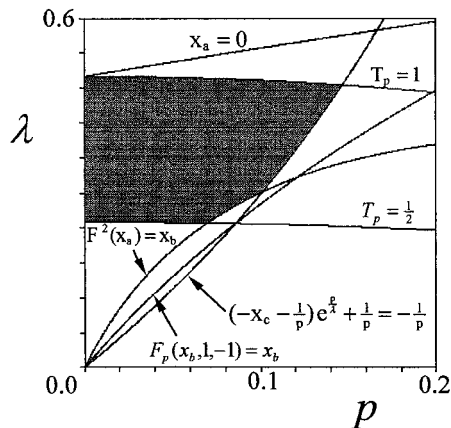


Fig. 6. Two-parameter diagram.

That is, the system has two symmetric periodic attractors, one of which is shown in Fig. 2(c).  $\square$

In this lemma, we can see an essential function of the ICC that makes stable dynamics by averaging two expanding maps with opposite slopes  $(d/d\bar{x})f(\bar{x}, 1) > 1$ ,  $(d/d\bar{x})f(\bar{x}, -1) < -1$ , and  $1/2|(d/d\bar{x})f(\bar{x}, 1) + (d/d\bar{x})f(\bar{x}, -1)| < 1$  for  $x_a < |\bar{x}| < x_b$ . Then Lemma 1 and Lemma 2 guarantee the coexisting phenomenon of chaos synchronization, a periodic attractor, and their symmetric ones. The parameters' conditions (8), (10) are satisfied in the shaded region in Fig. 6.

### III. CONCLUSION

Using the ICC, we have considered a simple coupling system of two nonautonomous chaotic circuits. The ICC changes the two chaotic attractors into a coexisting state of chaos synchronization, a periodic attractor and their symmetric ones. The coexisting phenomenon is guaranteed theoretically and is demonstrated in the laboratory. Now we are extending the ICC system to a coupled system of a large number of chaotic circuits and are analyzing their various synchronous phenomena. It may be developed into a novel artificial neural network.

### REFERENCES

- [1] H. Torikai and T. Saito, "Occasional linear connection for synchronization of chaos," *IEEE Trans. Circuits Syst. I*, vol. 43, pp. 374–385, May 1996.
- [2] —, "Synchronization of chaos and its itinerancy from a network by occasional linear connection," *IEEE Trans. Circuits Syst. I*, vol. 45, pp. 464–472, Apr. 1998.
- [3] J. J. Hopfield and A. V. M. Herz, "Rapid local synchronization of action potentials: Toward computation with coupled integrate-and-fire neurons," *Proc. Nat. Acad. Sci. USA*, vol. 92, pp. 6655–6662, 1995.
- [4] J. F. Heagy, L. M. Pecora, and T. L. Carroll, "Short wavelength bifurcations and size instabilities in coupled oscillator systems," *Phys. Rev. Lett.*, vol. 74, no. 21, pp. 4185–4188, 1995.
- [5] M. P. Kennedy, "Communication with chaos: State of the art and engineering challenges," in *Proc. NDES*, 1996, pp. 1–8.
- [6] H. Nijmeijer, I. Blekhan, A. Fradkov, and A. Pogromsky, "Self-synchronization and controlled synchronization," in *Proc. COC*, St. Petersburg, FL, 1997, pp. 36–41.
- [7] M. G. Rosenblum, A. S. Pikovsky, and J. Kurths, "Phase synchronization of chaotic oscillators," *Phys. Rev. Lett.*, vol. 76, pp. 1804–1807, 1996.
- [8] E. Ott and J. C. Sommerer, "Blowout bifurcation: The occurrence of riddled basins and on-off intermittency," *Phys. Lett., A*, vol. 188, pp. 39–47, 1994.
- [9] T. Stojanovski, L. Kocarev, and U. Parlitz, "Driving and synchronizing by chaotic impulses," *Phys. Rev. E*, vol. 54, no. 2, pp. 2128–2131, 1996.

- [10] T. Yang, L. B. Yang, and C. M. Yang, "Impulsive synchronization of Lorenz systems," *Phys. Lett. A*, vol. 226, pp. 349–354, 1997.
- [11] Y. Ohmori, M. Nakagawa, and T. Saito, "Mutual coupling of oscillators with chaos and period doubling bifurcation," in *Proc. IEEE/ISCAS*, 1986, pp. 61–64.
- [12] T. Y. Li and J. A. Yorke, "Ergodic transformation from an interval into itself," *Trans. Amer. Math. Soc.*, vol. 45, no. 5, pp. 464–472, 1998.

## Recurrent Least Squares Support Vector Machines

J. A. K. Suykens and J. Vandewalle

**Abstract**—The method of support vector machines (SVM's) has been developed for solving classification and static function approximation problems. In this paper we introduce SVM's within the context of recurrent neural networks. Instead of Vapnik's epsilon insensitive loss function, we consider a least squares version related to a cost function with equality constraints for a recurrent network. Essential features of SVM's remain, such as Mercer's condition and the fact that the output weights are a Lagrange multiplier weighted sum of the data points. The solution to recurrent least squares (LS-SVM's) is characterized by a set of nonlinear equations. Due to its high computational complexity, we focus on a limited case of assigning the squared error an infinitely large penalty factor with early stopping as a form of regularization. The effectiveness of the approach is demonstrated on trajectory learning of the double scroll attractor in Chua's circuit.

**Index Terms**—Double scroll, radial basis functions, recurrent neural networks, support vector machines.

### I. INTRODUCTION

Recently, support vector machines (SVM's) have been introduced as a new method for solving classification and function estimation problems with many successful applications [24]–[27]. SVM's are based on the structural risk minimization principle. The quality and complexity of the SVM solution does not depend directly on the dimensionality of the input space. The derivation of SVM's is based on constructing an optimal separating hyperplane after nonlinearly mapping the input space into a higher dimensional space. The explicit construction of this mapping is avoided by the application of Mercer's condition. Kernels that satisfy this condition and can be employed for SVM's are polynomials, splines, radial basis functions, and multilayer perceptrons with one hidden layer. For classification problems the parameters which are related to these kernel functions are chosen so as to minimize an upper bound on the Vapnik–Chervonenkis (VC) dimension of the SVM. The training of SVM's with Vapnik's epsilon insensitive loss function is done by quadratic programming. The number of hidden units in the

Manuscript received November 11, 1998; revised December 1, 1999 and January 14, 2000. This work was carried out at the ESAT Laboratory and the Interdisciplinary Center of Neural Networks ICNN of the Katholieke Universiteit Leuven, in the framework of the FWO Project G.0262.97 Learning and Optimization: an Interdisciplinary Approach, the Belgian Programme on Interuniversity Poles of Attraction, initiated by the Belgian State, Prime Minister's Office for Science, Technology, and Culture (IUAP P4-02 and IUAP P4-24) and the Concerted Action Project Model-Based Information Processing Systems (MIPS) of the Flemish Community. Johan Suykens is a Postdoctoral Researcher with the National Fund for Scientific research FWO-Flanders. This paper was recommended by Associate Editor J. M. Zurada.

The authors are with the Department of Electrical Engineering, ESAT-SISTA, Katholieke Universiteit Leuven, B-3001 Leuven (Heverlee), Belgium (e-mail: johan.suykens@esat.kuleuven.ac.be).

Publisher Item Identifier S 1057-7122(00)05504-5.

SVM is determined by the number of support vector data, which corresponds to the number of nonzero coefficients in the solution vector to the QP problem. A least squares version of SVM's for function estimation and classification has been investigated in [14] and [23]. The training is done then by solving a set of linear equations. On the other hand, sparseness is lost in the least squares version.

In this paper we introduce a recurrent neural network version of least squares SVM's (LS-SVM's). We investigate a class of nonlinear output error models for the autonomous case. In time-series prediction applications the difference between feedforward and recurrent is crucial, especially concerning chaotic systems [9], [20], [29]. Often the model is trained as a feedforward network (one step ahead predictor). On the other hand in order to make a further continuation of the given time series, the identified feedforward model has to be iterated, i.e., it has to be used as a recurrent network for prediction. Due to the sensitivity of chaotic systems for small perturbations, replacing the feedforward model by its recurrent version (replacing the true signal values by the estimated ones as input of the network) may result into large deviations. In this paper we consider the problem of trajectory learning of the double scroll attractor [4], [5], [10] by recurrent neural networks which are parametrized by SVM's. In [21] a simple recurrent neural state space model has been trained for the double scroll. It has been observed that by using classical dynamic backpropagation [11], [12] this is a difficult task [21].

For the recurrent LS-SVM's essential features of SVM's, such as Mercer's condition and the fact that the output weights are a Lagrange multiplier weighted sum of the data points, are still applicable. The solution to recurrent LS-SVM's is characterized by a set of nonlinear equations. Because of the considered least squares norm there seems to be some similarity at first sight with Widrow's adaline [9], [29]. However, recurrent LS-SVM's and Widrow's adaline are fundamentally different, e.g., due to the recurrent versus feedforward architecture, respectively, with corresponding implications for the training procedure.

The training problem for recurrent LS-SVM's is formulated as a nonconvex constrained nonlinear optimization problem in the error variables and the Lagrange multipliers. The computational cost of the training process is relatively expensive. Due to this high computational complexity, we focus on a limited case of assigning the squared error an infinitely large penalty factor with early stopping as a form of regularization. A sequential quadratic programming [7] training process is behaving well for this problem. For the example of trajectory learning of the double scroll attractor, recurrent LS-SVM's with a radial basis function (RBF) kernel yield a good generalization performance. Due to the study of the limited case, early stopping has to be applied which is equivalent to a form of regularization [1], [16]. The overfitting phenomenon can be detected on the training data itself by taking the first given data points in time as the initial state for the recurrent LS-SVM, while in the case of standard feedforward models, early stopping is usually based upon a test data set. Finally, part of the SVM theory that is related to upper bounds on the generalization error is not applicable because the input vectors to the recurrent architecture are not independent of each other.

This paper is organized as follows. In Section II we review work on classical SVM's and LS-SVM's. In Section III, recurrent LS-SVM's are introduced. In Section IV, recurrent LS-SVM's with an RBF kernel are applied to trajectory learning of the double scroll attractor.

## II. SUPPORT VECTOR MACHINES AND FUNCTION ESTIMATION

Here we review basic ideas of the support vector method of static function estimation and LS-SVM's. For detailed information about SVM's we refer to [3], [6], [15], [17]–[19], and [24]–[27].

Let us consider first the regression in the set of linear functions

$$\mathcal{F}(\mathcal{X}) = \mathcal{W}^T \mathcal{X} + \mathcal{B} \quad (1)$$

given training data  $\{\mathcal{X}_i, \mathcal{Y}_i\}_{i=1}^M$  where  $M$  denotes the number of training data,  $\mathcal{X}_i \in \mathbb{R}^m$  are the input data,  $\mathcal{Y}_i \in \mathbb{R}$  are the output data, and  $\mathcal{W} \in \mathbb{R}^m$ ,  $\mathcal{B} \in \mathbb{R}$ . Originally, in the support vector method one aims at minimizing the empirical risk

$$\mathcal{R}_{\text{emp}}(\mathcal{W}, \mathcal{B}) = \frac{1}{M} \sum_{i=1}^M |\mathcal{Y}_i - \mathcal{W}^T \mathcal{X}_i - \mathcal{B}|_\epsilon \quad (2)$$

subject to elements of a structure  $S_n$ , defined by the inequality  $\mathcal{W}^T \mathcal{W} \leq c_n$ . The loss function employs Vapnik's  $\epsilon$ -insensitive model

$$|\mathcal{Y}_i - \mathcal{F}(\mathcal{X}_i)|_\epsilon = \begin{cases} 0, & \text{if } |\mathcal{Y}_i - \mathcal{F}(\mathcal{X}_i)| \leq \epsilon \\ |\mathcal{Y}_i - \mathcal{F}(\mathcal{X}_i)| - \epsilon, & \text{otherwise.} \end{cases} \quad (3)$$

The function estimation problem is formulated then as

$$\min_{\mathcal{W}, \mathcal{B}, \xi^*, \xi} \mathcal{J}_\epsilon(\mathcal{W}, \xi^*, \xi) = \frac{1}{2} \mathcal{W}^T \mathcal{W} + \gamma \left\{ \sum_{i=1}^M \xi_i^* + \sum_{i=1}^M \xi_i \right\} \quad (4)$$

subject to the constraints

$$\begin{cases} \mathcal{Y}_i - \mathcal{W}^T \mathcal{X}_i - \mathcal{B} \leq \epsilon + \xi_i^*, & i = 1, \dots, M \\ -\mathcal{Y}_i + \mathcal{W}^T \mathcal{X}_i + \mathcal{B} \leq \epsilon + \xi_i, & i = 1, \dots, M \\ \xi_i^* \geq 0, & i = 1, \dots, M \\ \xi_i \geq 0, & i = 1, \dots, M \end{cases}$$

where  $\xi_i, \xi_i^*$  are slack variables and  $\gamma$  is a positive real constant. The solution is given by

$$\mathcal{W} = \sum_{i=1}^M (\alpha_i^* - \alpha_i) \mathcal{X}_i \quad (5)$$

where  $\alpha_i^*, \alpha_i$  are obtained by solving a quadratic program and are the Lagrange multipliers related to the first and second set of constraints. The data points corresponding to nonzero values for  $(\alpha_i^* - \alpha_i)$  in (5) are called support vectors. Typically, many of these values are equal to zero. Another loss function which has been investigated is

$$\mathcal{J}_{\epsilon, p}(\mathcal{W}, \xi^*, \xi) = \frac{1}{2} \mathcal{W}^T \mathcal{W} + \gamma \left\{ \sum_{i=1}^M (\xi_i^*)^p + \sum_{i=1}^M (\xi_i)^p \right\} \quad (6)$$

where  $p = 1$  corresponds to (4).

The work in this paper is related to a least squares version of SVM's. This version has been investigated in [14] for function estimation and in [23] for classification problems. For function estimation it corresponds to the following form of ridge regression

$$\min_{\mathcal{W}, \mathcal{B}, \xi} \mathcal{J}_{\text{LS}}(\mathcal{W}, \mathcal{B}, \xi) = \frac{1}{2} \mathcal{W}^T \mathcal{W} + \gamma \frac{1}{2} \sum_{i=1}^M \xi_i^2 \quad (7)$$

subject to the equality constraints

$$\mathcal{Y}_i = \mathcal{W}^T \mathcal{X}_i + \mathcal{B} + \xi_i, \quad i = 1, \dots, M.$$

One defines the Lagrangian

$$\begin{aligned} \mathcal{L}_{\text{LS}}(\mathcal{W}, \mathcal{B}, \xi; \alpha) &= \mathcal{J}_{\text{LS}}(\mathcal{W}, \mathcal{B}, \xi) \\ &\quad - \sum_{i=1}^M \alpha_i \left( \mathcal{W}^T \mathcal{X}_i + \mathcal{B} + \xi_i - \mathcal{Y}_i \right) \end{aligned} \quad (8)$$

where  $\alpha_i$  are Lagrange multipliers (which can be either positive or negative due to the equality constraints as follows from the Kuhn–Tucker conditions [7]). The conditions for optimality

$$\begin{cases} \frac{\partial \mathcal{L}_{LS}}{\partial \mathcal{W}} = 0 \rightarrow \mathcal{W} = \sum_{i=1}^M \alpha_i \mathcal{X}_i \\ \frac{\partial \mathcal{L}_{LS}}{\partial \mathcal{B}} = 0 \rightarrow \sum_{i=1}^M \alpha_i = 0 \\ \frac{\partial \mathcal{L}_{LS}}{\partial \xi_i} = 0 \rightarrow \alpha_i = \gamma \xi_i, & i = 1, \dots, M \\ \frac{\partial \mathcal{L}_{LS}}{\partial \alpha_i} = 0 \rightarrow \mathcal{W}^T \mathcal{X}_i + \mathcal{B} + \xi_i - \mathcal{Y}_i = 0, & i = 1, \dots, M \end{cases} \quad (9)$$

can be written immediately as the solution to the following set of linear equations:

$$\left[ \begin{array}{ccc|c} I & 0 & 0 & -\mathcal{X} \\ 0 & 0 & 0 & -\tilde{\mathbf{1}}^T \\ 0 & 0 & \gamma I & -I \\ \hline \mathcal{X}^T & \tilde{\mathbf{1}} & I & 0 \end{array} \right] \begin{bmatrix} \mathcal{W} \\ \mathcal{B} \\ \xi \\ \alpha \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad (10)$$

where  $\mathcal{X} = [\mathcal{X}_1 \dots \mathcal{X}_M]$ ,  $\mathcal{Y} = [\mathcal{Y}_1; \dots; \mathcal{Y}_M]$ ,  $\tilde{\mathbf{1}} = [1; \dots; 1]$ ,  $\xi = [\xi_1; \dots; \xi_M]$ ,  $\alpha = [\alpha_1; \dots; \alpha_M]$ . The solution is finally given by

$$\left[ \begin{array}{c|c} 0 & \tilde{\mathbf{1}}^T \\ \hline \tilde{\mathbf{1}} & \mathcal{X}^T \mathcal{X} + \gamma^{-1} I \end{array} \right] \begin{bmatrix} \mathcal{B} \\ \alpha \end{bmatrix} = \begin{bmatrix} 0 \\ \mathcal{Y} \end{bmatrix} \quad (11)$$

with  $\mathcal{W} = \sum_i \alpha_i \mathcal{X}_i$ ,  $\xi_i = \alpha_i / \gamma$ . The support values  $\alpha_i$  are proportional now to the errors at the data points (9).

So far we explained the linear case. SVM's with polynomials, splines, radial basis function networks, or multilayer perceptrons as kernels are obtained after mapping the input data into a higher dimensional space by  $\varphi(\mathcal{X}_i)$ , where  $\varphi(\cdot) : \mathbb{R}^m \rightarrow \mathbb{R}^{n_h}$ . The number  $n_h$  does not have to be specified because of the application of Mercer's condition, which means that

$$K(\mathcal{X}_i, \mathcal{X}_j) = \varphi(\mathcal{X}_i)^T \varphi(\mathcal{X}_j) \quad (12)$$

can be imposed for these kernels. For LS-SVM's the nonlinear function can then be represented as

$$\mathcal{Y}_i = \sum_{j=1}^M \alpha_j K(\mathcal{X}_j, \mathcal{X}_i) + \mathcal{B} \quad (13)$$

where the parameters  $\alpha_j, \mathcal{B}$  follow from (11) after replacing  $\mathcal{X}_i^T \mathcal{X}_j$  by  $K(\mathcal{X}_i, \mathcal{X}_j)$ . For RBF kernels one has

$$K(\mathcal{X}_i, \mathcal{X}_j) = \exp(-\nu \|\mathcal{X}_i - \mathcal{X}_j\|_2^2) \quad (14)$$

where  $\nu$  is a positive real constant. The number of hidden units in the LS-SVM is equal to the number of data points, while in standard SVM's this is equal to the number of support vectors. In other words, sparseness is lost in the LS-SVM case.

### III. RECURRENT LEAST SQUARES SVM'S

Given a deterministic nonlinear dynamical system with input  $u_k \in \mathbb{R}$  and output  $y_k \in \mathbb{R}$ , we consider nonlinear models of the form

$$\hat{y}_k = f(\hat{y}_{k-1}, \hat{y}_{k-2}, \dots, \hat{y}_{k-p}, u_{k-1}, u_{k-2}, \dots, u_{k-p}) \quad (15)$$

where  $\hat{y}_k$  denotes the estimated output and  $f$  is a smooth nonlinear mapping. Depending on the field, such models are called recurrent

input/output models [9], [20], [29], nonlinear output error (NOE) models [16], or parallel models [11]. These are opposed to

$$\hat{y}_k = f(y_{k-1}, y_{k-2}, \dots, y_{k-p}, u_{k-1}, u_{k-2}, \dots, u_{k-p}) \quad (16)$$

which are called feedforward models, NARX models or series-parallel models. Models of this form can be trained by means of the standard SVM methods discussed in Section II. The parametrization of  $f$  by SVM's is static because there is no recursion in the variable  $\hat{y}_k$ .

Without loss of generality concerning the methods, we will further discuss the autonomous case of the recurrent model (15)

$$\hat{y}_k = f(\hat{y}_{k-1}, \hat{y}_{k-2}, \dots, \hat{y}_{k-p}). \quad (17)$$

We take the following parametrization:

$$\hat{y}_k = w^T \varphi \left( \begin{bmatrix} \hat{y}_{k-1} \\ \hat{y}_{k-2} \\ \vdots \\ \hat{y}_{k-p} \end{bmatrix} \right) + b. \quad (18)$$

We express this model in terms of the given data and the error variables:

$$y_k - e_k = w^T \varphi(x_{k-1|k-p} - \xi_{k-1|k-p}) + b \quad (19)$$

where  $e_k = y_k - \hat{y}_k$ ,  $x_{k-1|k-p} = [y_{k-1}; y_{k-2}; \dots; y_{k-p}]$ ,  $\xi_{k-1|k-p} = [e_{k-1}; e_{k-2}; \dots; e_{k-p}]$  by definition. The nonlinear mapping  $\varphi(\cdot) : \mathbb{R}^p \rightarrow \mathbb{R}^{n_h}$  will be related to Mercer's condition. The output weight vector and bias term are denoted by  $w \in \mathbb{R}^{n_h}$  and  $b \in \mathbb{R}$ . Recurrent neural networks (15) and (18) are classically trained by dynamic backpropagation or backpropagation through time [11], [12], [28]. The goal of this Section is to develop an SVM approach for the recurrent model.

We formulate the training of the network (19) as

$$\min_{w, b, e} \mathcal{J}(w, b, e) = \frac{1}{2} w^T w + \gamma \frac{1}{2} \sum_{k=p+1}^{N+p} e_k^2 \quad (20)$$

subject to the equality constraints

$$y_k - e_k = w^T \varphi(x_{k-1|k-p} - \xi_{k-1|k-p}) + b, \quad k = p+1, \dots, N+p. \quad (21)$$

We define the Lagrangian

$$\begin{aligned} \mathcal{L}(w, b, e; \alpha) &= \mathcal{J}(w, b, e) + \sum_{k=p+1}^{N+p} \alpha_{k-p} \\ &\times \left[ y_k - e_k - w^T \varphi(x_{k-1|k-p} - \xi_{k-1|k-p}) - b \right]. \end{aligned} \quad (22)$$

The conditions for optimality are given by (23) at the bottom of the next page. By replacing  $w$  into the last two conditions and applying Mercer's condition by letting  $\mathcal{X}_i$  in (1) correspond to  $z_{k-1|k-p} = x_{k-1|k-p} - \xi_{k-1|k-p}$

$$K(z_{k-1|k-p}, z_{l-1|l-p}) = \varphi(z_{k-1|k-p})^T \varphi(z_{l-1|l-p}) \quad (24)$$

one obtains the following conditions for optimality in (25) at the bottom of the next page. Hence, (23) which is of the form  $F_1(w, b, e, \alpha) = 0$  has been represented as  $F_2(b, e, \alpha) = 0$  in (25) by elimination of  $w$ . This is similar to the elimination of  $w$  in (11) for the static SVM case

where  $w$  itself is never explicitly calculated. Furthermore, the Mercer condition (24) is implicitly defining the nonlinear mapping  $\varphi(\cdot)$ .

Finding a solution to (25) is computationally very expensive. Therefore, we further consider the case  $\gamma \rightarrow \infty$  which corresponds to

$$\min_{\alpha, e, b} \frac{1}{2} \sum_{k=p+1}^{N+p} e_k^2 \quad (26)$$

subject to (27) at the bottom of the page. At this point it is important to note that the vectors  $z$  depend on the error variables  $e$ . Hence, the cost function is subject to a set of nonlinear constraints in  $e$ . Hence, although (26) and (27) might seem similar to Widrow's adaline at first sight, due to the least squares norm, the recurrent LS-SVM is fundamentally different.

The resulting recurrent simulation model is given by

$$\hat{y}_k = \sum_{l=p+1}^{N+p} \alpha_{l-p} K \left( z_{l-1|l-p}, \begin{bmatrix} \hat{y}_{k-1} \\ \hat{y}_{k-2} \\ \vdots \\ \hat{y}_{k-p} \end{bmatrix} \right) + b \quad (28)$$

with given initial condition  $\hat{y}_i = y_i$  for  $i = 1, 2, \dots, p$ . For RBF kernels one employs

$$K(z_{k-1|k-p}, z_{l-1|l-p}) = \exp(-\nu \|z_{k-1|k-p} - z_{l-1|l-p}\|_2^2) \quad (29)$$

where  $\nu$  is a positive real constant.

No centers have to be determined in SVM models, in contrast with most of the classical RBF approaches [2], [13], [15]. For the recurrent SVM case the parameter estimation problem becomes nonconvex. The constrained nonlinear optimization problem (26), (27) can be solved, e.g., by sequential quadratic programming (SQP) [7]. For large data sets special methods for large scale nonlinear optimization have to be applied. Because the results are based on the limit case  $\gamma \rightarrow \infty$  and, hence, partially neglecting the regularization term  $(1/2)w^T w$  in (26) (the form of the solution is still derived by taking into account a regularization term based upon (20)) there is a danger for overfitting. Therefore, instead of minimizing the cost function to its local minimum one has to apply early stopping, which is equivalent to a form of regularization [1], [3], [16]. In standard SVM theory, which is applicable to feedforward models, the parameter  $\nu$  in (29) could be determined by minimizing upper bounds on the generalization error. However, this theory does not apply to the recurrent SVM case because the input arguments  $z$  of  $K(\cdot, \cdot)$  are not independent of each other. Hence,  $\nu$  should be either chosen as part of the unknown parameter vector for the optimization (26), (27) or ad hoc.

#### IV. EXAMPLE: DOUBLE SCROLL TRAJECTORY LEARNING

In this example we consider Chua's circuit [4], [5], [10]

$$\begin{cases} \dot{x} = \alpha[y - h(x)] \\ \dot{y} = x - y + z \\ \dot{z} = -\beta y \end{cases} \quad (30)$$

$$\begin{cases} \frac{\partial \mathcal{L}}{\partial w} = w - \sum_{k=p+1}^{N+p} \alpha_{k-p} \varphi(x_{k-1|k-p} - \xi_{k-1|k-p}) = 0 \\ \frac{\partial \mathcal{L}}{\partial b} = \sum_{k=p+1}^{N+p} \alpha_{k-p} = 0 \\ \frac{\partial \mathcal{L}}{\partial e_k} = \gamma e_k - \alpha_{k-p} - \sum_{i=1}^p \alpha_{k-p+i} \frac{\partial}{\partial e_{k-i}} [w^T \varphi(x_{k-1|k-p} - \xi_{k-1|k-p})] = 0, \quad k = p+1, \dots, N \\ \frac{\partial \mathcal{L}}{\partial \alpha_{k-p}} = y_k - e_k - w^T \varphi(x_{k-1|k-p} - \xi_{k-1|k-p}) - b = 0, \quad k = p+1, \dots, N+p. \end{cases} \quad (23)$$

$$\begin{cases} \sum_{k=p+1}^{N+p} \alpha_{k-p} = 0 \\ \gamma e_k - \alpha_{k-p} - \sum_{i=1}^p \alpha_{k-p+i} \frac{\partial}{\partial e_{k-i}} \left[ \sum_{l=p+1}^{N+p} \alpha_{l-p} K(z_{k-1|k-p}, z_{l-1|l-p}) \right] = 0, \quad k = p+1, \dots, N \\ y_k - e_k - \sum_{l=p+1}^{N+p} \alpha_{l-p} K(z_{k-1|k-p}, z_{l-1|l-p}) - b = 0, \quad k = p+1, \dots, N+p. \end{cases} \quad (25)$$

$$\begin{cases} y_k - e_k = \sum_{l=p+1}^{N+p} \alpha_{l-p} K(z_{l-1|l-p}, z_{k-1|k-p}) + b, \quad k = p+1, \dots, N+p \\ \sum_{k=p+1}^{N+p} \alpha_{k-p} = 0. \end{cases} \quad (27)$$

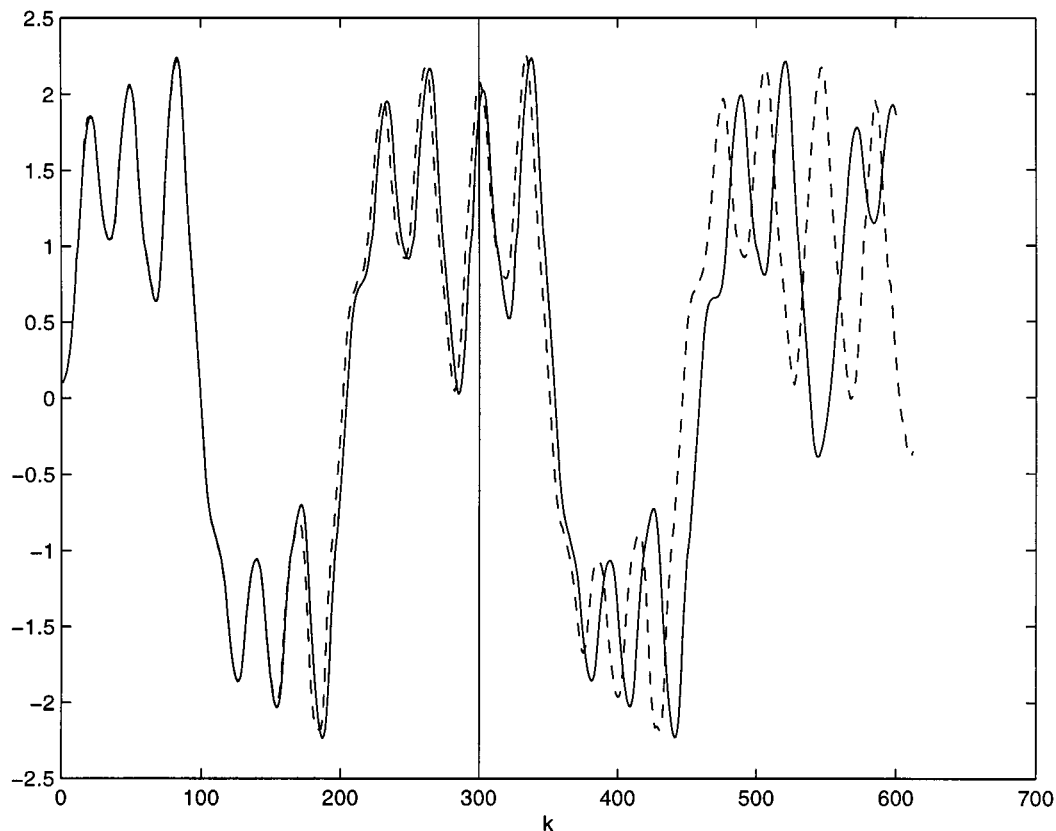


Fig. 1. Trajectory learning of the double scroll (full line) by an LS-SVM with an RBF kernel. The simulation result after training (dashed line) on  $N = 300$  data points is shown with as initial condition data points  $k = 1$  to 12. For the model structure  $p = 12$  and  $\nu = 1$  is taken. Early stopping is done in order to avoid overfitting.

with piecewise linear characteristic

$$h(x) = m_1 x + \frac{1}{2}(m_0 - m_1)(|x + 1| - |x - 1|). \quad (31)$$

A double scroll attractor is generated by taking  $\alpha = 9$ ,  $\beta = 14.286$ ,  $m_0 = -1/7$ ,  $m_1 = 2/7$ . A trajectory has been generated for initial condition  $[0.1; 0; -0.1]$  by using a Runge-Kutta integration rule (ode23 in Matlab).

In Fig. 1 the first  $N = 300$  data points were used for trajectory learning of the double scroll by a recurrent SVM with an RBF kernel. For the model structure  $p = 12$  has been taken and  $\nu = 1$  for the RBF kernel. In order to solve the constrained nonlinear optimization problem (26), (27), SQP has been applied (constr <AU: COMPLETE WORD? > in Matlab with a specification of the number of equality constraints). The model (28) has been simulated in C by using Matlab's cmex facility. In all simulations, for the initial unknown parameter vector  $\alpha_k, e_k$  have been chosen randomly according to a Gaussian distribution with zero mean and standard deviation 0.1 and  $b = 0$ . Although no further optimization of the model structure has been done, it has been observed that small values of  $p$  (that are chosen in accordance with Takens' embedding theorem) are slowing down the training process. The value of  $\nu$  in (29) has to be chosen relative to the scaling of the data set. In Fig. 1 the simulation result after training is shown in dashed line based on (28), (29) with as initial condition data points  $k = 1$  to 12. Early stopping has been decided based upon the given training data (not on an additional test set, as is usual for feedforward models) by evaluating the cost function in (26) for the simulation obtained from (28) with  $p = 12$  and  $N = 300$  for the given initial state. Early stopping was applied at the moment when this performance index degraded, which occurred after about 400 iteration steps of SQP for this

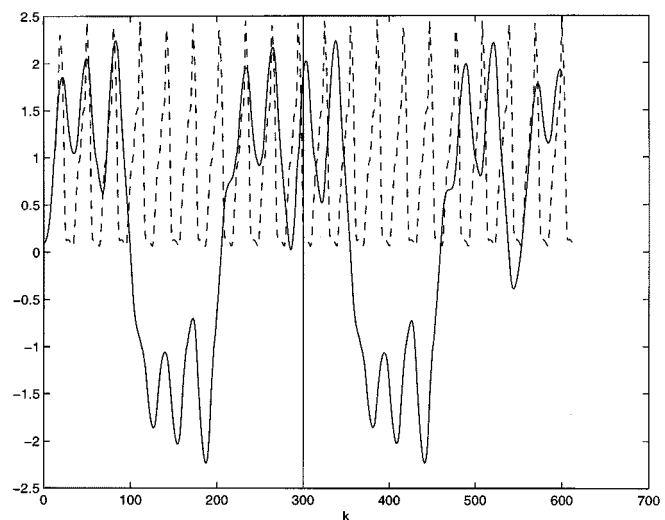


Fig. 2. Overfitting of the data obtained by taking the result of Fig. 1 as initial parameter vector for a further optimization to the local minimum.

example. The resulting recurrent SVM is generalizing to a double scroll attractor, also for small perturbations on the initial state. Fig. 2 shows a further optimization to the local minimum with the result of Fig. 1 as an initial parameter vector, which leads to overfitting. The overfitting phenomenon can be detected on the training data itself as illustrated on Fig. 2. Finally, note that for the recurrent SVM's it is important to simulate the form (28) because when solving the constrained nonlinear optimization problem (26), (27) the constraints will usually not hold exactly, but with a certain tolerance. For chaotic systems such small

perturbation could cause significant differences between the solution vector to (26), (27) and the recurrent simulation model (28).

## V. CONCLUSION

In this paper we introduced recurrent LS-SVM's, in addition to existing SVM solutions for classification and static nonlinear function estimation problems. Some essential features of SVM's such as Mercer's condition and the fact that the output weights are a Lagrange multiplier weighted sum of the data points, are still applicable. For computational reasons we focused on the limit case of assigning the squared error an infinitely large penalty factor. The training has been formulated as a nonconvex constrained nonlinear optimization problem in the error variables and the Lagrange multipliers. Training recurrent neural networks by dynamic backpropagation to follow trajectories of chaotic systems is well known to be a difficult problem. By trajectory learning of a double scroll we illustrated that recurrent LS-SVM's can generalize well, even on relatively small given training data sets. For large data sets, efficient methods for large scale constrained nonlinear optimization have to be used. The work of recurrent SVM's opens new perspectives with respect to time-series prediction and nonlinear modeling in general.

## REFERENCES

- [1] C. M. Bishop, *Neural Networks for Pattern Recognition*. New York: Oxford Univ. Press, 1995.
- [2] S. Chen, C. Cowan, and P. Grant, "Orthogonal least squares learning algorithm for radial basis function networks," *IEEE Trans. Neural Networks*, vol. 2, pp. 302–309, 1991.
- [3] V. Cherkassky and F. Mulier, *Learning from Data: Concepts, Theory and Methods*. New York: Wiley, 1998.
- [4] L. O. Chua, M. Komuro, and T. Matsumoto, "The double scroll family," *IEEE Trans. Circuits Syst. I*, vol. 33, pp. 1072–1118, Nov. 1986.
- [5] L. O. Chua, "Chua's circuit: An overview ten years later," *J. Circuits Syst. Comp.*, vol. 4, no. 2, pp. 117–159, 1994.
- [6] C. Cortes and V. Vapnik, "Support vector networks," *Mach. Learning*, vol. 20, pp. 273–297, 1995.
- [7] R. Fletcher, *Practical Methods of Optimization*. New York: Wiley, 1987.
- [8] G. H. Golub and C. F. Van Loan, *Matrix Computations*. Baltimore, MD: Johns Hopkins Univ. Press, 1989.
- [9] S. Haykin, *Neural Networks: A Comprehensive Foundation*. New York: Macmillan, 1994.
- [10] R. N. Madan, Ed., *Chua's Circuit: A Paradigm for Chaos*, Singapore: World Scientific, 1993.
- [11] K. S. Narendra and K. Parthasarathy, "Identification and control of dynamical systems using neural networks," *IEEE Trans. Neural Networks*, vol. 1, pp. 4–27, 1990.
- [12] —, "Gradient methods for the optimization of dynamical systems containing neural networks," *IEEE Trans. Neural Networks*, vol. 2, pp. 252–262, 1991.
- [13] T. Poggio and F. Girosi, "Networks for approximation and learning," *Proc. IEEE*, vol. 78, no. 9, pp. 1481–1497, 1990.
- [14] C. Saunders, A. Gammerman, and V. Vovk, "Ridge regression learning algorithm in dual variables," in *Proc. 15th Int. Conf. Machine Learning ICML-98*, Madison, WI, 1998.
- [15] B. Schölkopf, K.-K. Sung, C. Burges, F. Girosi, P. Niyogi, T. Poggio, and V. Vapnik, "Comparing support vector machines with Gaussian kernels to radial basis function classifiers," *IEEE Trans. Signal Processing*, vol. 45, pp. 2758–2765, Nov. 1997.
- [16] J. Sjöberg, Q. Zhang, L. Jung, A. Benveniste, B. Delyon, P. Glorennec, H. Hjalmarsson, and A. Juditsky, "Nonlinear black-box modeling in system identification: A unified overview," *Automatica*, vol. 31, no. 12, pp. 1691–1724, 1995.
- [17] A. Smola and B. Schölkopf, "On a kernel-based method for pattern recognition, regression, approximation and operator inversion," *Algoritmica*, vol. 22, pp. 211–231, 1998.
- [18] A. Smola, B. Schölkopf, and K.-R. Müller, "The connection between regularization operators and support vector kernels," *Neural Networks*, vol. 11, no. 4, 1998.
- [19] A. Smola, "Learning with kernels," Ph.D. thesis, GMD, Birlinghoven, 1999.
- [20] J. A. K. Suykens, J. Vandewalle, and B. De Moor, *Artificial Neural Networks for Modeling and Control of Non-Linear Systems*. Boston, MA: Kluwer, 1996.
- [21] J. A. K. Suykens and J. Vandewalle, "Learning a simple recurrent neural state space model to behave like Chua's double scroll," *IEEE Trans. Circuits Syst. I*, vol. 42, pp. 499–502, Aug. 1995.
- [22] —, *Nonlinear Modeling: Advanced Black-Box Techniques*. Boston, MA: Kluwer, 1998.
- [23] —, "Least squares support vector machine classifiers," *Neural Processing Lett.*, vol. 9, no. 3, pp. 293–300, June 1999.
- [24] V. Vapnik, *The Nature of Statistical Learning Theory*. New York: Springer-Verlag, 1995.
- [25] —, *Statistical Learning Theory*. New York: Wiley, 1998.
- [26] V. Vapnik, S. Golowich, and A. Smola, "Support vector method for function approximation, regression estimation and signal processing," in *Advances in Neural Information Processing Systems*. Cambridge, MA: MIT Press, 1997, vol. 9.
- [27] V. Vapnik, "The support vector method of function estimation," in *Nonlinear Modeling: Advanced Black-Box Techniques*, J. A. K. Suykens and J. Vandewalle, Eds. Boston, MA: Kluwer, 1998, pp. 55–85.
- [28] P. Werbos, "Backpropagation through time: What it does and how to do it," *Proc. IEEE*, vol. 78, no. 10, pp. 1150–1560, 1990.
- [29] J. M. Zurada, *Introduction to Artificial Neural Systems*. New York: West, 1992.