# SUPPORT VECTOR MACHINE FOR THE SIMULTANEOUS APPROXIMATION OF A FUNCTION AND ITS DERIVATIVE

M. Lázaro<sup>1</sup>, I. Santamaría<sup>2</sup>, F. Pérez-Cruz<sup>1</sup>, A. Artés-Rodríguez<sup>1</sup>

 <sup>1</sup>Departamento de Teoría de la Señal y Comunicaciones Universidad Carlos III, Leganés, 28911, Madrid, Spain e-mail: {marce,fernandop,antonio}@ieee.org
 <sup>2</sup>Departamento de Ingeniería de Comunicaciones Universidad de Cantabria, 39005 Santander, Spain e-mail: nacho@gtas.dicom.unican.es

Abstract. In this paper, the problem of simultaneously approximating a function and its derivative is formulated within the Support Vector Machine (SVM) framework. The problem has been solved by using the  $\varepsilon$ -insensitive loss function and introducing new linear constraints in the approximation of the derivative. The resulting quadratic problem can be solved by Quadratic Programming (QP) techniques. Moreover, a computationally efficient Iterative Re-Weighted Least Square (IRWLS) procedure has been derived to solve the problem in large data sets. The performance of the method has been compared with the conventional SVM for regression, providing outstanding results.

# INTRODUCTION

Regression approximation of a given data set is a very common problem in a number of applications. In some of these applications, like economy, device modeling, telemetry, etc, it is necessary to fit not only the underlying characteristic function but also its derivatives, which are often available. Some methods have been employed to simultaneously approximate a set of samples of a function and its derivative: splines, neural networks or filter bank-based methods are some examples (see [1] and references therein).

On the other hand, Support Vector Machines are state-of-the-art tools for linear and nonlinear input-output knowledge discovery [5, 4]. The Support Vector Machines, given a labeled dataset  $((\mathbf{x}_i, y_i))$ , where  $\mathbf{x}_i \in \mathbb{R}^d$  for  $i = 1, \ldots, N$ ) and a function  $\phi(\cdot)$  that nonlinearly transforms the input vector  $\mathbf{x}_i$  to a higher dimensional space, solve either classification  $(y_i \in \{\pm 1\})$  or regression  $(y_i \in \mathbb{R})$  problems. In this paper, we will deal with the regression approximation problem and we will extend the SVM framework when prior knowledge about the derivative of the functional relation between  $\mathbf{x}$  and y is known. We will solve this problem for one-dimensional problem (d = 1), but it can be readily extended to multidimensional input and the gradient information, as we will show herein. We would like to find the functional relation between x and ygiving a labeled data set  $((x_i, y_i, y'_i))$ , where  $y_i \in \mathbb{R}$  and  $y'_i \in \mathbb{R}$ , which is the derivative of the function to be approximated at  $x_i$ ).

We will solve it using the  $\varepsilon$ -insensitive loss function and it will lead to a solution similar to the SVM in which we have Support Vectors related to the function value  $(y_i)$  and Support Vectors related to the derivative value  $(y'_i)$ , and both together form the complete SVM expansion for regression approximation with information about the derivative of the function.

The solution to the proposed algorithm is obtained using an Iterative Re-Weighted Least Square (IRWLS) procedure that it is able to easily solve the SVM. This procedure has been successfully applied to the regular SVM for classification [2] for regression [3].

The rest of the paper is outlined as follows. The modification of the SVM for regression approximation to include the derivative information is presented in Section 2. In Section 3, we deal with the IRWLS procedure for linear machines. The use of nonlinear kernel and the application of the kernel trick to the IRWLS procedure is presented in Section 4. The generalization of the method to higher order input spaces is outlined in Section 5. Experimental results are presented in Section 6. We conclude the paper in Section 7 with some final comments and possible further works.

#### PROPOSED SVM-BASED APPROACH

As it has been outlined before, the proposed method is an extension of the Support Vector Machine for Regression (SVM-R) employing the Vapnik's  $\varepsilon$ insensitive loss function [5]. The SVM-R obtains a linear regressor in the transformed space (feature space)

$$f(x) = \mathbf{w}^T \phi(x) + b, \tag{1}$$

where  $\mathbf{w}$  and b define the linear regression, which is nonlinear in the input space (unless  $\phi(x)$  is linear). Roughly speaking, the SVM-R minimizes the squared norm of the weight vector  $\mathbf{w}$ , while linearly penalizes the deviations greater than  $\varepsilon$ .

With respect to the conventional SVM-R cost function, the proposed method adds a new penalty term: the errors in the derivative that are out of its associated insensitive region. In the general case, a different parameter is employed to define the insensitive region size for the function ( $\varepsilon$ ) and for the derivative ( $\varepsilon'$ ). Taking this extension into account, the proposed approach minimizes

$$J(\mathbf{w}, \boldsymbol{\xi}, \boldsymbol{\xi}^*, \boldsymbol{\tau}, \boldsymbol{\tau}^*) = \frac{1}{2} ||\mathbf{w}||^2 + C_1 \sum_{i=1}^N (\xi_i + \xi_i^*) + C_2 \sum_{i=1}^N (\tau_i + \tau_i^*), \quad (2)$$

subject to

$$\mathbf{w}^T \phi(x_i) + b - y_i \le \varepsilon + \xi_i, \tag{3}$$

$$y_i - \mathbf{w}^T \phi(x_i) - b \le \varepsilon + \xi_i^*, \tag{4}$$

$$\mathbf{w}^T \phi'(x_i) - y_i' \le \varepsilon' + \tau_i, \tag{5}$$

$$y_i' - \mathbf{w}^T \phi'(x_i) \le \varepsilon' + \tau_i^*, \tag{6}$$

$$\xi_i, \xi_i^*, \tau_i, \tau_i^* \ge 0. \tag{7}$$

The positive slack variables  $\xi$ ,  $\xi^*$ ,  $\tau$  and  $\tau^*$  are responsible for penalizing errors greater than  $\varepsilon$  and  $\varepsilon'$ , respectively, in the function and derivative, and  $\phi'(x)$  denotes the derivative of  $\phi(x)$ . To solve this problem, the following Lagrangian functional is employed, introducing the previous linear constraints

$$L(\mathbf{w}, \boldsymbol{\xi}, \boldsymbol{\xi}^{*}, \boldsymbol{\tau}, \boldsymbol{\tau}^{*}, \boldsymbol{\alpha}, \boldsymbol{\alpha}^{*}, \boldsymbol{\lambda}, \boldsymbol{\lambda}^{*}) = \frac{1}{2} ||\mathbf{w}||^{2} + C_{1} \sum_{i=1}^{N} (\xi_{i} + \xi_{i}^{*}) + C_{2} \sum_{i=1}^{N} (\tau_{i} + \tau_{i}^{*})$$
$$- \sum_{i=1}^{N} \alpha_{i} [\varepsilon + \xi_{i} - (\mathbf{w}^{T} \phi(x_{i}) + b - y_{i})] - \sum_{i=1}^{N} \alpha_{i}^{*} [\varepsilon + \xi_{i}^{*} - (y_{i} - \mathbf{w}^{T} \phi(x_{i}) - b)]$$
$$- \sum_{i=1}^{N} \lambda_{i} [\varepsilon' + \tau_{i} - (\mathbf{w}^{T} \phi'(x_{i}) - y'_{i})] - \sum_{i=1}^{N} \lambda_{i}^{*} [\varepsilon' + \tau_{i}^{*} - (y'_{i} - \mathbf{w}^{T} \phi'(x_{i}))]$$
$$- \sum_{i=1}^{N} (\mu_{i} \xi_{i} + \mu_{i}^{*} \xi_{i}^{*} + \theta_{i} \tau_{i} + \theta_{i}^{*} \tau_{i}^{*}). \quad (8)$$

This functional has to be minimized with respect to  $\mathbf{w}$ , b,  $\xi$ ,  $\xi^*$ ,  $\tau$  and  $\tau^*$ , and maximized with respect to the Lagrange multipliers. The solution to this problem can be obtained considering the Karush-Kuhn-Tucker (KKT) complementary conditions, which for this specific problem are

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{w} + \sum_{i=1}^{N} \alpha_i \phi(x_i) - \sum_{i=1}^{N} \alpha_i^* \phi(x_i) + \sum_{i=1}^{N} \lambda_i \phi'(x_i) - \sum_{i=1}^{N} \lambda_i^* \phi'(x_i) = 0, \quad (9)$$

$$\frac{\partial L}{\partial b} = \sum_{i=1}^{N} (\alpha_i - \alpha_i^*) = 0, \qquad (10)$$

$$\frac{\partial L}{\partial \xi_i} = C_1 - \alpha_i - \mu_i = 0, \qquad \frac{\partial L}{\partial \xi_i^*} = C_1 - \alpha_i^* - \mu_i^* = 0, \qquad (11)$$

$$\frac{\partial L}{\partial \tau_i} = C_2 - \lambda_i - \theta_i = 0, \qquad \frac{\partial L}{\partial \tau_i^*} = C_2 - \lambda_i^* - \theta_i^* = 0, \qquad (12)$$

$$\alpha_i[\varepsilon + \xi_i - (\mathbf{w}^T \phi(x_i) + b - y_i)] = 0,$$
(13)

$$\alpha_i^* [\varepsilon + \xi_i^* - (y_i - \mathbf{w}^T \phi(x_i) - b)] = 0, \qquad (14)$$

$$\Lambda_i[\varepsilon' + \tau_i - (\mathbf{w}^T \phi'(x_i) - y_i')] = 0, \qquad (15)$$

$$\Lambda_{i}^{*}[\varepsilon' + \tau_{i}^{*} - (y_{i}' - \mathbf{w}^{T}\phi'(x_{i}))] = 0, \qquad (16)$$

and

$$\mu_i \xi_i = 0, \quad \mu_i^* \xi_i^* = 0, \quad \theta_i \tau_i = 0, \quad \theta_i^* \tau_i^* = 0.$$
(17)

From (9), the weight vector  $\mathbf{w}$  takes the form

$$\mathbf{w} = \sum_{i=1}^{N} (\alpha_i^* - \alpha_i) \phi(x_i) + \sum_{i=1}^{N} (\lambda_i^* - \lambda_i) \phi'(x_i),$$
(18)

which means that the regression is

$$f(x) = \sum_{i=1}^{N} (\alpha_i^* - \alpha_i) < \phi(x_i), \phi(x) > + \sum_{i=1}^{N} (\lambda_i^* - \lambda_i) < \phi'(x_i), \phi(x) > +b.$$
(19)

Substituting (18) into (8), rearranging terms and taking into account (10)-(12), one thus arrives to the Wolfe's dual problem needing to maximize

$$W(\boldsymbol{\alpha}, \boldsymbol{\alpha}^{*}, \boldsymbol{\lambda}, \boldsymbol{\lambda}^{*}) = -\frac{1}{2} \sum_{i,j=1}^{N} (\alpha_{i}^{*} - \alpha_{i})(\alpha_{j}^{*} - \alpha_{j}) < \phi(x_{i}), \phi(x_{j}) >$$

$$-\frac{1}{2} \sum_{i,j=1}^{N} (\lambda_{i}^{*} - \lambda_{i})(\lambda_{j}^{*} - \lambda_{j}) < \phi'(x_{i}), \phi'(x_{j}) >$$

$$-\frac{1}{2} \sum_{i,j=1}^{N} (\alpha_{i}^{*} - \alpha_{i})(\lambda_{j}^{*} - \lambda_{j}) < \phi(x_{i}), \phi'(x_{j}) >$$

$$-\frac{1}{2} \sum_{i,j=1}^{N} (\lambda_{i}^{*} - \lambda_{i})(\alpha_{j}^{*} - \alpha_{j}) < \phi'(x_{i}), \phi(x_{j}) >$$

$$-\sum_{i=1}^{N} [(\alpha_{i} + \alpha_{i}^{*})\varepsilon + (\lambda_{i} + \lambda_{i}^{*})\varepsilon'] + \sum_{i=1}^{N} [(\alpha_{i}^{*} - \alpha_{i})y_{i} + (\lambda_{i}^{*} - \lambda_{i})y_{i}'], \quad (20)$$

subject to

$$0 \le \alpha_i, \alpha_i^* \le C_1, \tag{21}$$

$$0 \le \lambda_i, \lambda_i^* \le C_2, \tag{22}$$

and

$$\sum_{i=1}^{N} (\alpha_i - \alpha_i^*) = 0.$$
(23)

It can be seen that (20) is a quadratic functional that only depends on the Lagrange multipliers  $\alpha_i$ ,  $\alpha_i^*$ ,  $\lambda_i$  and  $\lambda_i^*$ . This problem can be solved by Quadratic Programming (QP) techniques. Moreover, in the SVM framework, the nonlinear transformation  $\phi(x)$  is not needed to be explicitly known and it can be replaced by its kernels. In this case,  $\langle \phi(x_i), \phi(x_j) \rangle$  is substituted by  $K(x_i, x_j)$ , a kernel satisfying the Mercer Theorem [4]. From this definition for the kernel, it is easy to demonstrate that,

$$\langle \phi'(x_i), \phi(x_j) \rangle = \frac{\partial K(x_i, x_j)}{\partial x_i} \triangleq K'(x_i, x_j),$$
 (24)

$$\langle \phi(x_i), \phi'(x_j) \rangle = \frac{\partial K(x_i, x_j)}{\partial x_j} \triangleq G(x_i, x_j),$$
 (25)

and

$$\langle \phi'(x_i), \phi'(x_j) \rangle = \frac{\partial^2 K(x_i, x_j)}{\partial x_i \partial x_j} \triangleq J(x_i, x_j).$$
 (26)

Though K must be a Mercer Kernel, its derivatives do not have to. Therefore, using a valid kernel  $K(x_i, x_j)$ , once the Lagrange multipliers have been obtained, the regression estimate takes the form

$$f(x) = \sum_{i=1}^{N} (\alpha_i^* - \alpha_i) K(x_i, x) + \sum_{i=1}^{N} (\lambda_i^* - \lambda_i) K'(x_i, x) + b.$$
(27)

## **IRWLS ALGORITHM**

The QP solution of the system can be computationally expensive, especially when a large number of samples is employed. In this case the computational burden can make the problem unaffordable. In order to reduce the computational burden, an Iterative Re-Weighted Least Square (IRWLS) procedure has been developed. This IRWLS algorithm follows the same basic idea proposed in [3]. First at all, the Lagrangian (8) is rearranged to group the term depending on  $\xi_i$ ,  $\xi_i^*$ ,  $\tau_i$  and  $\tau_i^*$ . Taking into account (11)-(12), these terms can be eliminated. Therefore, (8) can be written as

$$L = \frac{1}{2} ||\mathbf{w}||^2 - \sum_{i=1}^{N} \alpha_i [\varepsilon + y_i - \mathbf{w}^T \phi(x_i) - b] - \sum_{i=1}^{N} \alpha_i^* [\varepsilon + \mathbf{w}^T \phi(x_i) + b - y_i] - \sum_{i=1}^{N} \lambda_i [\varepsilon' + y_i' - \mathbf{w}^T \phi'(x_i)] - \sum_{i=1}^{N} \lambda_i^* [\varepsilon' + \mathbf{w}^T \phi'(x_i) - y_i'].$$
(28)

This functional can be rewritten as

$$L = \frac{1}{2} ||\mathbf{w}||^2 + \frac{1}{2} \sum_{i=1}^{N} (a_i e_i^2 + a_i^* (e_i^*)^2) + \frac{1}{2} \sum_{i=1}^{N} (s_i d_i^2 + s_i^* (d_i^*)^2), \quad (29)$$

where

$$e_i = \mathbf{w}^T \phi(x_i) + b - y_i - \varepsilon, \quad a_i = \frac{2\alpha_i}{e_i}$$
 (30)

$$e_i^* = y_i - \mathbf{w}^T \phi(x_i) - b - \varepsilon, \quad a_i^* = \frac{2\alpha_i^*}{e_i^*}$$
(31)

$$d_i = \mathbf{w}^T \phi'(x_i) - y'_i - \varepsilon', \quad s_i = \frac{2\lambda_i}{d_i}$$
(32)

$$d_i^* = y_i' - \mathbf{w}^T \phi'(x_i) - \varepsilon', \quad s_i^* = \frac{2\lambda_i^*}{d_i^*}$$
(33)

This functional can be seen as a weighted least square one, where  $e_i$ , and  $e_i^*$  are the prediction error over the function,  $d_i$  and  $d_i^*$  are the prediction error over the derivative, and  $a_i$ ,  $a_i^*$ ,  $s_i$  and  $s_i^*$  are the corresponding weights. It is necessary to iterate because  $a_i = a_i(e_i)$ ,  $a_i^* = a_i^*(e_i^*)$ ,  $s_i = s_i(d_i)$  and  $s_i^* = s_i^*(d_i^*)$ . The goal at each iteration is to minimize (28) with respect to  $\mathbf{w}$  and b, supposing that  $a_i$ ,  $a_i^*$ ,  $s_i$  and  $s_i^*$  are fixed. Taking the derivative with respect to both variables, the following two equations are obtained:

$$\begin{bmatrix} \boldsymbol{\Phi}^{T} \mathbf{D}_{\mathbf{a}+\mathbf{a}^{*}} \boldsymbol{\Phi} + \boldsymbol{\Phi}^{\prime T} \mathbf{D}_{\mathbf{s}+\mathbf{s}^{*}} \boldsymbol{\Phi}^{\prime} + \mathbf{I} & \boldsymbol{\Phi}^{T} \mathbf{D}_{\mathbf{a}+\mathbf{a}^{*}} \mathbf{1} \\ (\mathbf{a}+\mathbf{a}^{*})^{T} \boldsymbol{\Phi} & (\mathbf{a}+\mathbf{a}^{*})^{T} \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix} = \\ \begin{bmatrix} \boldsymbol{\Phi}^{T} [\mathbf{D}_{\mathbf{a}+\mathbf{a}^{*}} \mathbf{y} + \mathbf{D}_{\mathbf{a}-\mathbf{a}^{*}} \mathbf{1}\varepsilon] + \boldsymbol{\Phi}^{\prime T} [\mathbf{D}_{\mathbf{s}+\mathbf{s}^{*}} \mathbf{y}^{\prime} + \mathbf{D}_{\mathbf{s}-\mathbf{s}^{*}} \mathbf{1}\varepsilon^{\prime}] \\ (\mathbf{a}+\mathbf{a}^{*})^{T} \mathbf{y} + (\mathbf{a}-\mathbf{a}^{*})^{T} \mathbf{1}\varepsilon \end{bmatrix}, \quad (34)$$

where  $\mathbf{a}, \mathbf{a}^*, \mathbf{s}$  and  $\mathbf{s}^*$  are the vectors containing the N corresponding weights in (29),  $\mathbf{D}_{\mathbf{a}}$  denotes the diagonal matrix with vector  $\mathbf{a}$  in the diagonal, and

$$\Phi = [\phi(x_1), \phi(x_2), \cdots, \phi(x_N)]^T 
\Phi' = [\phi'(x_1), \phi'(x_2), \cdots, \phi'(x_N)]^T$$
(35)

## **IRWLS WITH KERNELS**

The system (34) can be used when  $\phi(x)$  is known. When working with kernels, it is necessary to obtain the Lagrange multipliers to provide the regression (27). In this case, it must be taken into account that the weight vector can be written as

$$\mathbf{w} = [\mathbf{\Phi}^T, \mathbf{\Phi}'^T] \begin{bmatrix} \boldsymbol{\beta} \\ \boldsymbol{\gamma} \end{bmatrix}, \tag{36}$$

where  $\beta = \alpha^* - \alpha$  and  $\gamma = \lambda^* - \lambda$ . Now (36) will be used to replace **w** in (34). The first equation can be written as

$$[\boldsymbol{\Phi}^{T}\mathbf{D}_{\mathbf{a}+\mathbf{a}^{*}}\boldsymbol{\Phi}+\boldsymbol{\Phi}^{\prime T}\mathbf{D}_{\mathbf{s}+\mathbf{s}^{*}}\boldsymbol{\Phi}^{\prime}+\mathbf{I}]\mathbf{w} = [\boldsymbol{\Phi}^{T},\boldsymbol{\Phi}^{\prime T}] \begin{bmatrix} \mathbf{D}_{\mathbf{a}+\mathbf{a}^{*}}(\mathbf{y}-\mathbf{1}b) + \mathbf{D}_{\mathbf{a}-\mathbf{a}^{*}}\mathbf{1}\varepsilon \\ \mathbf{D}_{\mathbf{s}+\mathbf{s}^{*}}\mathbf{y}^{\prime} + \mathbf{D}_{\mathbf{s}-\mathbf{s}^{*}}\mathbf{1}\varepsilon^{\prime} \end{bmatrix}$$
(37)

Substituting  $\mathbf{w}$  by (36), multiplying in both sides by

$$\left( \left[ \boldsymbol{\Phi}^{T}, \boldsymbol{\Phi}^{\prime T} \right]^{T} \right)^{+} = \left( \left[ \begin{array}{c} \boldsymbol{\Phi} \\ \boldsymbol{\Phi}^{\prime} \end{array} \right] \left[ \boldsymbol{\Phi}^{T}, \boldsymbol{\Phi}^{\prime T} \right] \right)^{-1} \left[ \begin{array}{c} \boldsymbol{\Phi} \\ \boldsymbol{\Phi}^{\prime} \end{array} \right], \tag{38}$$

making mathematical arrangements and taking into account that we have two sets of decoupled equations, this equation can be written down as

$$\mathbf{H}^{-1} \left\{ \mathbf{H} \begin{bmatrix} \mathbf{D}_{\mathbf{a}+\mathbf{a}^{*}} & 0\\ 0 & \mathbf{D}_{\mathbf{s}+\mathbf{s}^{*}} \end{bmatrix} \mathbf{H} + \mathbf{H} \right\} \begin{bmatrix} \boldsymbol{\beta} \\ \boldsymbol{\gamma} \end{bmatrix} = \\ \begin{bmatrix} \mathbf{D}_{\mathbf{a}+\mathbf{a}^{*}}(\mathbf{y}-\mathbf{1}b) + \mathbf{D}_{\mathbf{a}-\mathbf{a}^{*}}\mathbf{1}\varepsilon \\ \mathbf{D}_{\mathbf{s}+\mathbf{s}^{*}}\mathbf{y}' + \mathbf{D}_{\mathbf{s}-\mathbf{s}^{*}}\mathbf{1}\varepsilon' \end{bmatrix} .$$
(39)

**H** can be easily computed from the definition of  $K(x_i, x_j)$  and (24)-(26)

$$\mathbf{H} = \begin{bmatrix} \boldsymbol{\Phi} \boldsymbol{\Phi}^T & \boldsymbol{\Phi} \boldsymbol{\Phi}'^T \\ \boldsymbol{\Phi}' \boldsymbol{\Phi}^T & \boldsymbol{\Phi}' \boldsymbol{\Phi}'^T \end{bmatrix} = \begin{bmatrix} \mathbf{K} & \mathbf{G} \\ \mathbf{K}' & \mathbf{J} \end{bmatrix},$$
(40)

where  $\mathbf{K}|_{ij} = K(x_i, x_j)$ ,  $\mathbf{K}'|_{ij} = K'(x_i, x_j)$ ,  $\mathbf{G}|_{ij} = G(x_i, x_j)$  and  $\mathbf{J}|_{ij} = J(x_i, x_j)$ . Canceling  $\mathbf{H}^{-1}\mathbf{H}$ , multiplying by the inverse of the diagonal matrix in (39) and moving b to the first term, the equation can be simplified. Finally, instead of using the second equation in (34), the simpler constraint (10) can be used, leading to the following whole system

$$\begin{bmatrix} \mathbf{H} + \begin{bmatrix} \mathbf{D}_{\mathbf{a}+\mathbf{a}^*} & 0 \\ 0 & \mathbf{D}_{\mathbf{s}+\mathbf{s}^*} \\ \begin{bmatrix} \mathbf{1}^T, \mathbf{0}^T \end{bmatrix}^{-1} & \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} \boldsymbol{\beta} \\ \boldsymbol{\gamma} \\ b \end{bmatrix} = \begin{bmatrix} \mathbf{y} + \frac{\mathbf{a}-\mathbf{a}^*}{\mathbf{a}+\mathbf{a}^*}\varepsilon \\ \mathbf{y}' + \frac{\mathbf{s}-\mathbf{s}^*}{\mathbf{s}+\mathbf{s}^*}\varepsilon' \\ 0 \end{bmatrix}, \quad (41)$$

where  $\frac{\mathbf{a}-\mathbf{a}^*}{\mathbf{a}+\mathbf{a}^*}$  denotes the vector containing  $(a_i - a_i^*)/(a_i + a_i^*)$  in the *i*-th row.

# Recalculating the weights

After each iteration, the new values for the weights  $a_i$ ,  $a_i^*$ ,  $s_i$  and  $s_i^*$  have to be calculated. It is simple to demonstrate, looking at the KKT conditions, that the corresponding values are given by

$$a_{i} = \begin{cases} 0, & e_{i} < 0 \\ \frac{2C_{1}}{e_{i}}, & e_{i} > 0 \end{cases}, \qquad a_{i}^{*} = \begin{cases} 0, & e_{i}^{*} < 0 \\ \frac{2C_{1}}{e_{i}^{*}}, & e_{i}^{*} > 0 \end{cases}$$
(42)

$$s_{i} = \begin{cases} 0, \quad d_{i} < 0 \\ \frac{2C_{2}}{d_{i}}, \quad d_{i} > 0 \end{cases}, \qquad s_{i}^{*} = \begin{cases} 0, \quad d_{i}^{*} < 0 \\ \frac{2C_{2}}{d_{i}^{*}}, \quad d_{i}^{*} > 0 \end{cases}$$
(43)

It must be noted that in (41) the inverse of these weights is used. In any case, when the estimated errors are null or very low, a maximum limit has to be imposed to the weight constants. This has the effect of adding a small number in the diagonal of the system in this case. This guarantees the matrix system is invertible. With this numerical trick to ensure convergence, the IRWLS algorithm can be summarized in Table 1.

1. Initialization:

- Compute H (from K, K', G and J)
- $a_i = C_1, s_i = C_2$  for odd  $i; a_i^* = C_1, s_i^* = C_2$  for even i.
- 2. To solve (41)
- 3. To evaluate

$$\mathbf{e} = \mathbf{K}^T \boldsymbol{\beta} + \mathbf{K}'^T \boldsymbol{\gamma} + \mathbf{1}b - \mathbf{y} - \mathbf{1}\varepsilon, \quad \mathbf{e}^* = \mathbf{y} - \mathbf{K}^T \boldsymbol{\beta} - \mathbf{K}'^T \boldsymbol{\gamma} - \mathbf{1}b - \mathbf{1}\varepsilon \\ \mathbf{d} = \mathbf{G}^T \boldsymbol{\beta} + \mathbf{J}^T \boldsymbol{\gamma} - \mathbf{y}' - \mathbf{1}\varepsilon', \qquad \mathbf{d}^* = \mathbf{y}' - \mathbf{G}^T \boldsymbol{\beta} - \mathbf{J}^T \boldsymbol{\gamma} - \mathbf{1}\varepsilon'$$

- 4. Recalculate  $a_i, a_i^*, s_i^*$  and  $s_i^*$  by (42) and (43) (with a maximum limit)
- 5. Go to step 2 until convergence is achieved

TABLE 1: IRWLS ALGORITHM PSEUDOCODE

## **EXTENSION OF THE METHOD**

The proposed method can be easily extended to d-dimensional input spaces and to consider up to k-th order derivatives. It is only necessary to incorporate the corresponding constraints. Because of the space limitation we have omitted the development, but in this case the solution takes the form

$$\mathbf{w} = \sum_{i=1}^{N} \sum_{j_1=0}^{k} \cdots \sum_{j_d=0}^{k} (\lambda_{ij_1\cdots j_d}^* - \lambda_{ij_1\cdots j_d}) \frac{\partial^{(j_1+\cdots+j_d)}\phi(\mathbf{x}_i)}{\partial x_1^{j_1}, \cdots, x_d^{j_d}},$$
(44)

where  $\mathbf{x}_i = [x_{i1}, x_{i2}, \cdots, x_{id}]^T$  and  $\lambda^*_{ij_1\cdots j_d}$  and  $\lambda_{ij_1\cdots j_d}$  are the Lagrange multipliers associated to the constraint in the *i*-th sample of

$$\frac{\partial^{(j_1+\dots+j_d)}f(\mathbf{x})}{\partial x_1^{j_1},\dots,x_d^{j_d}}$$

## RESULTS

In this section, some experimental results show the advantages of this method in the reconstruction of the derivative with respect to the conventional SVM-R approach. As test functions, we have selected a set of bandlimited functions: specifically, in each experiment a linear combination of 100 sinusoids with random amplitudes, frequencies (between 0 and 1 Hz) and phases has been generated. In the first example, 100 equally spaced sampling points in the range 0-5 have been employed by the SVM-R (100 samples of the function), And 50 points by the proposed approach (50 samples of the function +50 samples of the derivative). In this way, the number of total available data is the same. 1000 experiments have been considered, with a signal to noise ratio (SNR) of 20 dB in the samples of both the function and the derivative. Figure 1 plots the mean values of signal to error ratio (SER) in the reconstruction of the function (a) and of the derivative (b) as a function of the insensitivity parameter  $\varepsilon$ . In this case,  $\varepsilon' = \pi \varepsilon$  has been considered to take into account the different amplitude range of function and derivative (the mean amplitude of the derivative is  $\pi$  times higher than the mean amplitude of the function).



It can be seen that the proposed method (labeled SVM-D) provides better results than the SVM-R, specially in the reconstruction of the derivative. Moreover, a similar number of support vectors has been observed for both methods in all simulations (in the proposed method: support vectors related to the function + support vectors related to the derivative). Therefore, the proposed method does not increment the storage requirements of the model. The number of support vectors, as a function of  $\varepsilon$ , is ploted in Figure 2.



Figure 2: Number of total support vectors as a function of  $\varepsilon$ 

Finally, this method reduces the sensitivity of the method to the selection of  $\sigma$  for the Gaussian kernels. Figure 3 shows the SER in the reconstruction of the function as a function of  $\sigma$  for  $\varepsilon = 0.5$ . It can be seen that the  $\sigma$  range to obtain high SER values is increased by using the proposed method.



#### CONCLUSIONS

A new SVM-based method for the simultaneous reconstruction of a function and its derivative has been presented. A computationally efficient IRWLS algorithm has been derived to allow the application of the method to large data sets. This method provides better results than the conventional SVM-R approach in the reconstruction of function and derivative even when the same number of labeled data is employed in both methods. In this case, the proposed method needs a similar number of support vectors than conventional SVM-R. Moreover, the inclusion of the information of the derivatives reduces the dependence on the kernel size for Gaussian kernels. Expression have been presented for a one-dimensional input space and only the first derivative, but the extension of the method to higher order input spaces and derivative orders is straightforward by adding the corresponding constraints to the regression. Further work in necessary to simplify the parameter selection similarly to the  $\nu$ -SVM [4].

#### REFERENCES

- M. Lázaro, I. Santamaría, C. Pantaleón, J. Ibáñez and L. Vielva, "A Regularized technique for the simultaneous reconstruction of a function and its derivatives with application to nonlinear transistor modeling," Signal Processing. In Press, 2003.
- [2] F. Pérez-Cruz, A. Navia-Vázquez, J. L. Rojo-Álvarez and A. Artés-Rodríguez, "A New Training Algorithm for Support Vector Machines," in Proceedings of the Fifth Bayona Workshop on Emerging Technologies in Telecommunications, Baiona, Spain, Sept. 1999, pp. 116–120.
- [3] F. Pérez-Cruz, A. Navia-Vázquez, P. Alarcón-Diana and A. Artés-Rodríguez, "An IRWLS procedure for SVR," in Proceedings of the EUSIPCO'00, Tampere, Finland, Sept. 2000.
- [4] B. Schölkopf and A. Smola, Learning with Kernels, MIT Press, 2002.
- [5] V. N. Vapnik, Statistical Learning Theory, New York: Wiley & Sons, 1998.