



ELSEVIER

Neurocomputing 48 (2002) 85–105

NEUROCOMPUTING

www.elsevier.com/locate/neucom

## Weighted least squares support vector machines: robustness and sparse approximation<sup>☆</sup>

J.A.K. Suykens\*, J. De Brabanter, L. Lukas, J. Vandewalle

*Department of Electrical Engineering, Katholieke Universiteit Leuven, ESAT-SISTA Kasteelpark Arenberg 10, B-3001 Leuven (Heverlee), Belgium*

Received 3 November 2000; accepted 6 June 2001

### Abstract

Least squares support vector machines (LS-SVM) is an SVM version which involves equality instead of inequality constraints and works with a least squares cost function. In this way, the solution follows from a linear Karush–Kuhn–Tucker system instead of a quadratic programming problem. However, sparseness is lost in the LS-SVM case and the estimation of the support values is only optimal in the case of a Gaussian distribution of the error variables. In this paper, we discuss a method which can overcome these two drawbacks. We show how to obtain robust estimates for regression by applying a weighted version of LS-SVM. We also discuss a sparse approximation procedure for weighted and unweighted LS-SVM. It is basically a pruning method which is able to do pruning based upon the physical meaning of the sorted support values, while pruning procedures for classical multilayer perceptrons require the computation of a Hessian matrix or its inverse. The methods of this paper are illustrated for RBF kernels and demonstrate how to obtain robust estimates with selection of an appropriate number of hidden units, in the case of outliers or non-Gaussian error distributions with heavy tails. © 2002 Elsevier Science B.V. All rights reserved.

**Keywords:** Support vector machines; (Weighted) least squares; Ridge regression; Sparse approximation; Robust estimation

<sup>☆</sup> This research 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 Learning and Optimization: an Interdisciplinary Approach, the Belgian Program on Interuniversity Poles of Attraction, initiated by the Belgian State, Prime Minister's Office for Science, Technology and Culture (IUAP P4-02 & IUAP P4-24) and the Concerted Action Project MEFISTO of the Flemish Community. Johan Suykens is a postdoctoral researcher with the National Fund for Scientific Research FWO - Flanders.

\* Corresponding author. Tel.: +32-16-32-18-02; fax: +32-16-32-19-70.

E-mail address: johan.suykens@esat.kuleuven.ac.be (J.A.K. Suykens).

## 1. Introduction

Support vector machines (SVM) for classification and nonlinear function estimation, as introduced by Vapnik [35,36] and further investigated by many others [4,21–25], is an important new methodology in the area of neural networks and nonlinear modelling [27]. While classical neural networks approaches, such as multilayer perceptrons (MLP) and radial basis function (RBF) networks [2,17], suffer from problems like the existence of many local minima and the choice of the number of hidden units [2,13], SVM solutions are characterized by convex optimization problems, up to the determination of a few additional tuning parameters. Moreover, model complexity follows from this convex optimization problem. Typically, one solves a convex quadratic programming (QP) problem in dual space in order to determine the SVM model. The formulation of the optimization problem in the primal space associated with this QP problem involves inequality constraints. In the case of function estimation it is related to Vapnik's epsilon-insensitive loss function. An interesting property of the SVM solution is that one obtains a sparse approximation, in the sense that many elements in the QP solution vector are equal to zero. The additional hyperparameters of the SVM model are often determined by model selection based upon generalization bounds, which have been derived within the area of statistical learning theory. SVM is a kernel based approach, which allows the use of linear, polynomial and RBF kernels and others that satisfy Mercer's condition.

Recently, least squares (LS) versions of SVM's have been investigated for classification [28] and function estimation [20]. In these LS-SVM formulations one works with equality instead of inequality constraints and a sum squared error (SSE) cost function as it is frequently used in training of classical neural networks. This reformulation greatly simplifies the problem in such a way that the solution is characterized by a linear system, more precisely a KKT (Karush–Kuhn–Tucker) system [8], which takes a similar form as the linear system that one solves in every iteration step by interior point methods for standard SVM's [26]. This linear system can be efficiently solved by iterative methods such as conjugate gradient [29]. However, despite these computationally attractive features, LS-SVM solutions also have some potential drawbacks. The first drawback is that sparseness is lost in the LS-SVM solution. In this case every data point is contributing to the model and the relative importance of a data point is given by its support value. The second drawback is that it is well known that the use of a SSE cost function without regularization might lead to estimates which are less robust, e.g. with respect to outliers on the data or when the underlying assumption of a Gaussian distribution for the error variables is not realistic.

The aim of this paper is to show that one can overcome these drawbacks concerning sparseness and robustness within the present LS-SVM framework. First of all one should note that only the output weights of the SVM model follow as solution to the linear system. Only these parameters are related to the SSE

cost function which means that (up to a certain extent) one can still correct for wrong assumptions by an appropriate choice of the hyperparameters. These can be determined in several possible ways such as crossvalidation, bootstrapping, VC bounds, Bayesian inference etc. [4,34]. We show in this paper how one can apply weighted least squares in order to produce a more robust estimate. This is done by first applying an (unweighted) LS-SVM and, in the second stage, associate weighting values to the error variables based upon the resulting error variables from the first stage. Techniques of weighted least squares are well known e.g. in statistics, identification and control theory and signal processing. For LS-SVM's it can be employed as a cheap and efficient way to make the solution robust. In this way it can be used as an alternative to other methods in robust estimation [14] as  $L_1$  estimators and  $M$ -estimators with Huber loss function. Robust estimation is also possible within the standard SVM context [36]. In standard SVM methodology, one chooses a given cost function (any convex cost function can be taken in principle as shown in [26]). Instead of this top-down approach the weighted LS-SVM approach aims at working bottom-up by solving a sequence of weighted LS-SVM's starting from the unweighted version. In this way one implicitly tries to find an optimal underlying cost function, instead of imposing the cost function beforehand. In this sense there is also a close link between solving the SVM problem for a given convex cost function by interior point methods and iterative weighting of LS-SVM solutions.

Furthermore, in this paper we illustrate how sparseness can be imposed to the weighted LS-SVM solution by gradually pruning the sorted support value spectrum. While in pruning methods for MLP's [2] (like optimal brain damage [16] and optimal brain surgeon [12]) the procedure involves a Hessian matrix or its inverse, the pruning in LS-SVM's can be done based upon the solution vector itself (note that this implicitly requires inverting the system). Less meaningful data points as indicated by their support values, are removed and the LS-SVM is re-computed on the basis of the remaining points while validating on the complete training data set. In this paper we focus on a cheap and simple pruning method. Other methods for obtaining a sparse approximation with LS-SVM's are possible [6]. The advantage of the method shown in this paper is that the determination of the hyperparameters can be kept localized, while in the other approaches one needs to solve the convex optimization problem (which implicitly corresponds to solving a sequence of linear systems) for a given set of hyperparameters.

In general, parametric models like MLP's or RBF networks are applicable within a broad range of applications of either static problems (classification, regression, density estimation) or dynamic problems (recurrent networks, optimal control). Standard SVM's on the other hand have only been applied to static problems. However, the use of equality constrained and SSE based formulations within LS-SVM greatly simplifies the formulations and allows us to extend the method to recurrent networks [31] and control applications [32]. In the latter case convexity is lost,

but some of the other interesting SVM properties remain applicable. The results of this paper on weighted LS-SVM's for robust estimation and sparse approximation further motivates the LS-SVM approach towards its use as a new and general neural network methodology.

This paper is organized as follows. In Section 2 we review some basic notions of LS-SVM's for function estimation. In Section 3 we discuss the weighted LS-SVM formulation. The sparse approximation procedure is discussed in Section 4. Finally, illustrative examples are given in Section 5.

## 2. LS-SVM for nonlinear function estimation

Given a training data set of  $N$  points  $\{x_k, y_k\}_{k=1}^N$  with input data  $x_k \in \mathbb{R}^n$  and output data  $y_k \in \mathbb{R}$ , one considers the following optimization problem in primal weight space:

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

such that

$$y_k = w^T \varphi(x_k) + b + e_k, \quad k = 1, \dots, N$$

with  $\varphi(\cdot): \mathbb{R}^n \rightarrow \mathbb{R}^{n_h}$  a function which maps the input space into a so-called higher dimensional (possibly infinite dimensional) feature space, weight vector  $w \in \mathbb{R}^{n_h}$  in primal weight space, error variables  $e_k \in \mathbb{R}$  and bias term  $b$ . Note that the cost function  $J$  consists of a *SSE* fitting error and a regularization term, which is also a standard procedure for the training of MLP's and is related to ridge regression [10]. The relative importance of these terms is determined by the positive real constant  $\gamma$ . In the case of noisy data one avoids overfitting by taking a smaller  $\gamma$  value. SVM problem formulations of this form have been investigated independently in [20] (without bias term) and [28].

In primal weight space one has the model

$$y(x) = w^T \varphi(x) + b. \quad (2)$$

The weight vector  $w$  can be infinite dimensional, which makes a calculation of  $w$  from (1) impossible in general. Therefore, one computes the model in the dual space instead of the primal space. One defines the Lagrangian

$$\mathcal{L}(w, b, e; \alpha) = J(w, e) - \sum_{k=1}^N \alpha_k \{w^T \varphi(x_k) + b + e_k - y_k\} \quad (3)$$

with Lagrange multipliers  $\alpha_k \in \mathbb{R}$  (called support values). The conditions for optimality are given by

$$\begin{cases} \frac{\partial \mathcal{L}}{\partial w} = 0 \rightarrow w = \sum_{k=1}^N \alpha_k \varphi(x_k), \\ \frac{\partial \mathcal{L}}{\partial b} = 0 \rightarrow \sum_{k=1}^N \alpha_k = 0, \\ \frac{\partial \mathcal{L}}{\partial e_k} = 0 \rightarrow \alpha_k = \gamma e_k, & k = 1, \dots, N, \\ \frac{\partial \mathcal{L}}{\partial \alpha_k} = 0 \rightarrow w^T \varphi(x_k) + b + e_k - y_k = 0, & k = 1, \dots, N. \end{cases} \quad (4)$$

These conditions are similar to standard SVM optimality conditions, except for the condition  $\alpha_k = \gamma e_k$ . At this point one loses the sparseness property in LS-SVM's [9].

After elimination of  $w, e$  one obtains the solution

$$\left[ \begin{array}{c|c} 0 & 1_v^T \\ \hline 1_v & \Omega + \frac{1}{\gamma} I \end{array} \right] \begin{bmatrix} b \\ \alpha \end{bmatrix} = \begin{bmatrix} 0 \\ y \end{bmatrix} \quad (5)$$

with  $y = [y_1; \dots; y_N]$ ,  $1_v = [1; \dots; 1]$ ,  $\alpha = [\alpha_1; \dots; \alpha_N]$  and  $\Omega_{kl} = \varphi(x_k)^T \varphi(x_l)$  for  $k, l = 1, \dots, N$ . According to Mercer's condition, there exists a mapping  $\varphi$  and an expansion

$$K(x, y) = \sum_i \varphi_i(x) \varphi_i(y), \quad x, y \in \mathbb{R}^n, \quad (6)$$

if and only if, for any  $g(x)$  such that  $\int g(x)^2 dx$  is finite, one has

$$\int K(x, y) g(x) g(y) dx dy \geq 0. \quad (7)$$

As a result, one can choose a kernel  $K(\cdot, \cdot)$  such that

$$K(x_k, x_l) = \varphi(x_k)^T \varphi(x_l), \quad k, l = 1, \dots, N. \quad (8)$$

The resulting LS-SVM model for function estimation becomes

$$y(x) = \sum_{k=1}^N \alpha_k K(x, x_k) + b \quad (9)$$

where  $\alpha, b$  are the solution to (5). We focus on the choice of an RBF kernel  $K(x_k, x_l) = \exp\{-\|x_k - x_l\|_2^2 / \sigma^2\}$  for the sequel.

When solving large linear systems, it is necessary to apply iterative methods [10] to (5). However, the matrix in (5) is not positive definite. According to [29] one can transform the system into a positive definite system such that iterative methods (conjugate gradient, successive overrelaxation or others) can be applied to it. Note that the computational complexity of the conjugate gradient method for solving a linear system  $\mathcal{A}x = \mathcal{B}$  is  $\mathcal{O}(Nr^2)$  where  $\text{rank}(\mathcal{C}) = r$  with  $\mathcal{A} = I + \mathcal{C}$  and  $\mathcal{A} \in \mathbb{R}^{N \times N}$ . It converges in at most  $r + 1$  steps. The speed of convergence depends on the condition number of the matrix. In the case of LS-SVM this is influenced by the choice of  $(\gamma, \sigma)$  when using an RBF kernel.

### 3. Robust estimation by weighted LS-SVM

In order to obtain a robust estimate based upon the previous LS-SVM solution, in a subsequent step, one can weight the error variables  $e_k = \alpha_k/\gamma$  by weighting factors  $v_k$ . This leads to the optimization problem:

$$\min_{w^\star, b^\star, e^\star} J(w^\star, e^\star) = \frac{1}{2} w^{\star T} w^\star + \frac{1}{2} \gamma \sum_{k=1}^N v_k e_k^{\star 2} \quad (10)$$

such that

$$y_k = w^{\star T} \varphi(x_k) + b^\star + e_k^\star, \quad k = 1, \dots, N.$$

The Lagrangian becomes

$$\mathcal{L}(w^\star, b^\star, e^\star; \alpha^\star) = J(w^\star, e^\star) - \sum_{k=1}^N \alpha_k^\star \{w^{\star T} \varphi(x_k) + b^\star + e_k^\star - y_k\}. \quad (11)$$

The unknown variables for this weighted LS-SVM problem are denoted by the  $\star$  symbol. From the conditions for optimality and elimination of  $w^\star, e^\star$  one obtains the KKT system

$$\left[ \begin{array}{c|c} 0 & 1_v^T \\ \hline 1_v & \Omega + V_\gamma \end{array} \right] \begin{bmatrix} b^\star \\ \alpha^\star \end{bmatrix} = \begin{bmatrix} 0 \\ y \end{bmatrix}, \quad (12)$$

where the diagonal matrix  $V_\gamma$  is given by

$$V_\gamma = \text{diag} \left\{ \frac{1}{\gamma v_1}, \dots, \frac{1}{\gamma v_N} \right\}. \quad (13)$$

The choice of the weights  $v_k$  is determined based upon the error variables  $e_k = \alpha_k/\gamma$  from the (unweighted) LS-SVM case (5). Robust estimates are obtained then (see

[5,18]) e.g. by taking

$$v_k = \begin{cases} 1 & \text{if } |e_k/\hat{s}| \leq c_1, \\ \frac{c_2 - |e_k/\hat{s}|}{c_2 - c_1} & \text{if } c_1 \leq |e_k/\hat{s}| \leq c_2, \\ 10^{-4} & \text{otherwise,} \end{cases} \quad (14)$$

where  $\hat{s}$  is a robust estimate of the standard deviation of the LS-SVM error variables  $e_k$ :

$$\hat{s} = \frac{\text{IQR}}{2 \times 0.6745}. \quad (15)$$

The interquartile range IQR is the difference between the 75th percentile and 25th percentile. In the estimate  $\hat{s}$  one takes into account how much the estimated error distribution deviates from a Gaussian distribution. Another robust estimate of the standard deviation is  $\hat{s} = 1.483 \text{ MAD}(x_i)$  where MAD stands for the median absolute deviation [11]. The SSE cost function in the unweighted LS-SVM formulation is optimal under the assumption of a normal Gaussian distribution for  $e_k$ . The procedure (14) corrects for this assumption in order to obtain a robust estimate when this distribution is not normal. Eventually, the procedure (10) (14) can be repeated iteratively, but in practice one single additional weighted LS-SVM step will often be sufficient. One assumes that  $e_k$  has a symmetric distribution which is usually the case when  $(\gamma, \sigma)$  are well-determined by an appropriate model selection method. The constants  $c_1, c_2$  are typically chosen as  $c_1 = 2.5$  and  $c_2 = 3$  [18]. This is a reasonable choice taking into account the fact that for a Gaussian distribution, there will be very few residuals larger than  $2.5\hat{s}$ . Another possibility is to determine  $c_1, c_2$  from a density estimation of the  $e_k$  distribution. Using these weightings one can correct for  $y$ -outliers or for a non-Gaussian instead of Gaussian error distributions.

This leads us to the following algorithm:

#### Algorithm 1—weighted LS-SVM

1. Given training data  $\{x_k, y_k\}_{k=1}^N$ , find an optimal  $(\gamma, \sigma)$  combination (e.g. by 10-fold cross-validation or generalization bounds) by solving linear systems (5). For the optimal  $(\gamma, \sigma)$  combination one computes  $e_k = \alpha_k/\gamma$  from (5).
2. Compute  $\hat{s}$  from the  $e_k$  distribution.
3. Determine the weights  $v_k$  based upon  $e_k, \hat{s}$ .
4. Solve the Weighted LS-SVM (12), giving the model  $y(x) = \sum_{k=1}^N \alpha_k^\star K(x, x_k) + b^\star$ .

An important notion in robust estimation is the breakdown point of an estimator [1,6,19]. Loosely speaking, it is the smallest fraction of contamination of a given data set that can result in an estimate which is arbitrarily far away from the

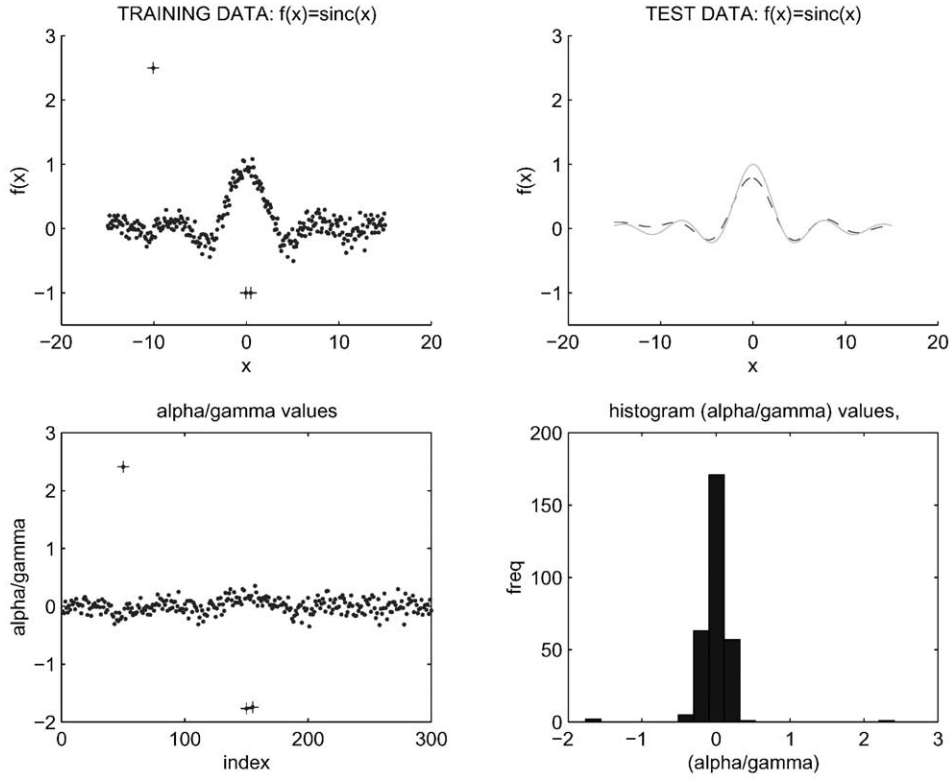


Fig. 1. Estimation of a sinc function by LS-SVM with RBF kernel, given 300 training data point, corrupted by zero mean Gaussian noise and three outliers (denoted by '+'). (Top-Left) Training data set; (Top-Right) resulting LS-SVM model evaluated on an independent test set: (solid line) true function, (dashed line) LS-SVM estimate; (Bottom-Left)  $e_k = \alpha_k/\gamma$  values; (Bottom-Right) histogram for distribution of  $e_k$  values, which is non-Gaussian due to the three outliers.

estimated parameter vector obtained from the uncontaminated data set. It is well known that the least squares estimate in linear regression (parametric) without regularization has a low breakdown point. At this point the unweighted LS-SVM has already better properties due to the fact that least squares is only applied for finding  $\alpha$  (which corresponds to estimating the output layer and not  $(\gamma, \sigma)$ ). It is still possible then to correct with  $(\gamma, \sigma)$  when using an RBF kernel. Although unweighted LS-SVM already has rather desirable properties, the breakdown point is further improved by applying the weighted LS-SVM afterwards.

In standard SVM approaches one usually chooses a certain cost function (any convex cost function can be taken according to [26]). However, in such a top-down approach one assumes in fact that the noise distribution is known because one knows which cost function is optimal for a given noise distribution. The weighted



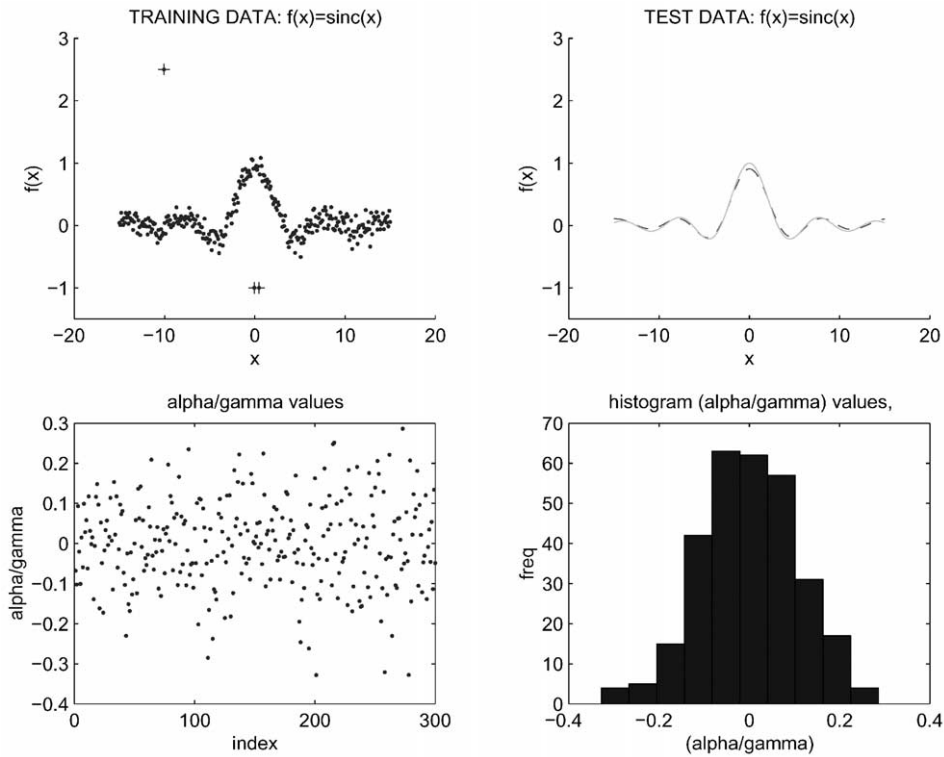


Fig. 2. (Continued) Weighted LS-SVM applied to the results of Fig. 1. The  $e_k$  distribution becomes Gaussian and the generalization performance on the test data improves.

LS-SVM procedure tries to work bottom–up and aims at finding the optimal underlying cost function for the given noisy data. The methods can be conceptually linked by the fact that when one applies an interior point method for solving a standard SVM problem, at every iteration step a KKT system is solved which is similar to solving one single LS-SVM. Hence, solving a standard SVM with arbitrary convex cost function implicitly corresponds, in fact, to solving a sequence of LS-SVMs.

#### 4. Imposing sparseness

While standard SVM's possess a sparseness property in the sense that many  $\alpha_k$  values are equal to zero, this is not the case for LS-SVM's due to the fact that  $\alpha_k = \gamma e_k$  from the conditions for optimality. An equivalence has also been proven between sparse approximation and SVM's [3,9]. Here we discuss a simple procedure that shows how a sparse approximation for LS-SVM's can be obtained. This

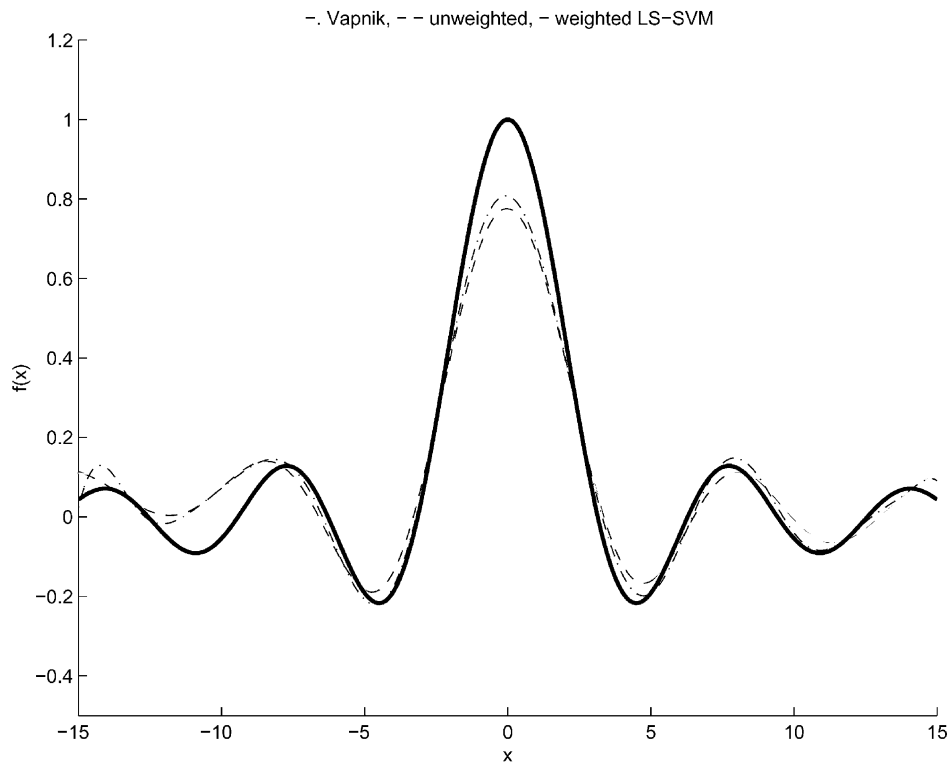


Fig. 3. (Continued) Comparison between standard SVM, unweighted LS-SVM and weighted LS-SVM. Weighted LS-SVM gives the best estimate for this example.

is done by exploiting the fact that the support values have a physical meaning in the sense that they reveal the relative importance of the data points for contributing to the model.

Here, we propose a pruning procedure for LS-SVM's which is based upon the sorted support value spectrum. By omitting a relative and small amount of the least meaningful data points (this corresponds to setting these  $\alpha_k$  values to zero) and re-estimating the LS-SVM, one obtains a sparse approximation. In order to guarantee a good generalization performance, in each of these pruning steps one can optimize  $(\gamma, \sigma)$  (by defining an independent validation set, 10-fold cross-validation, generalization bounds or others). An important difference with pruning methods for MLP's [2,12,16], e.g. optimal brain damage and optimal brain surgeon, is that pruning can be done based upon solution vector itself and does not require the knowledge of a Hessian matrix or its inverse. This pruning method has been successfully applied to unweighted LS-SVM's [30]. Here, we outline the algorithm for the weighted LS-SVM case.

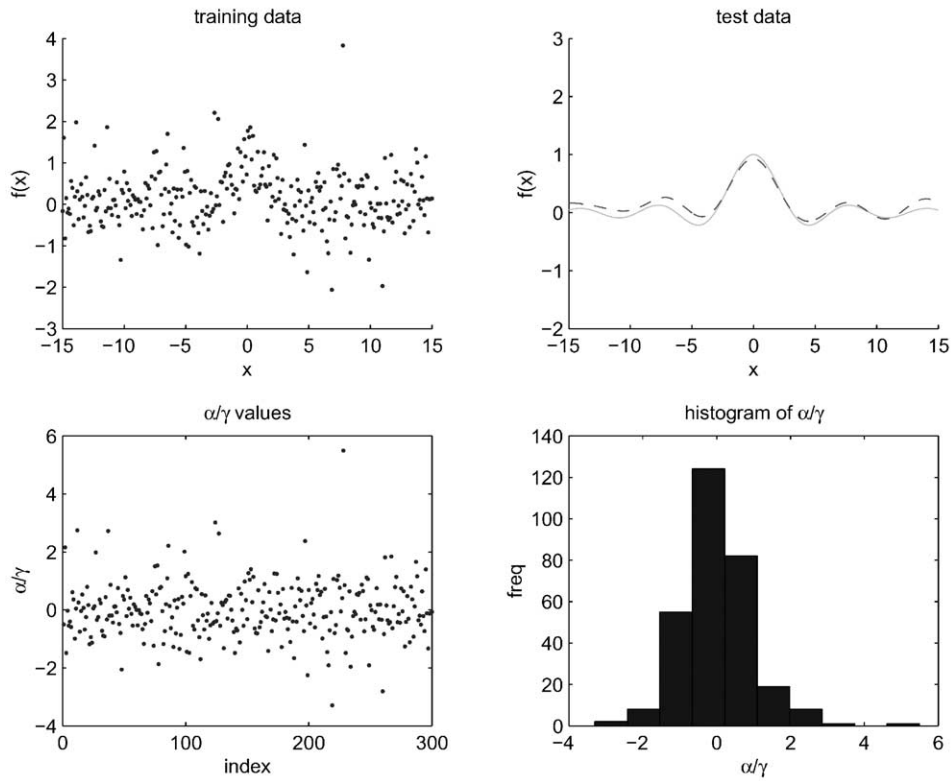


Fig. 4. Estimation of a sinc function by LS-SVM with RBF kernel, given 300 training data point, corrupted by a central  $t$ -distribution with heavy tails. (Top-Left) Training data set; (Top-Right) resulting LS-SVM model evaluated on an independent test set: (solid line) true function, (dashed line) LS-SVM estimate; (Bottom-Left)  $e_k = \alpha_k/\gamma$  values; (Bottom-Right) histogram for distribution of  $e_k$  values, which is non-Gaussian.

#### Algorithm 2—weighted LS-SVM pruning

1. Set  $N = N_{tot}$  equal to the number of training data.
2. Given  $N$  training data, apply Algorithm 1 where  $(\gamma, \sigma)$  is determined on the total amount of  $N_{tot}$  training data. The solution to the linear system for optimal  $(\gamma, \sigma)$  yields  $\alpha_k^\star$ .
3. Sort the values  $|\alpha_k^\star|$ .
4. Remove a small amount of  $M$  points (typically 5% of the  $N$  points) that have the smallest values in the sorted  $|\alpha_k^\star|$  spectrum.
5. Retain  $N - M$  points and set  $N := N - M$ .
6. Go to 2 and retrain on the reduced training set, unless the user-defined performance index degrades.

Typically, doing a number of pruning steps without modification of  $(\gamma, \sigma)$  will be possible. When the generalization performance starts degrading (checked e.g. on

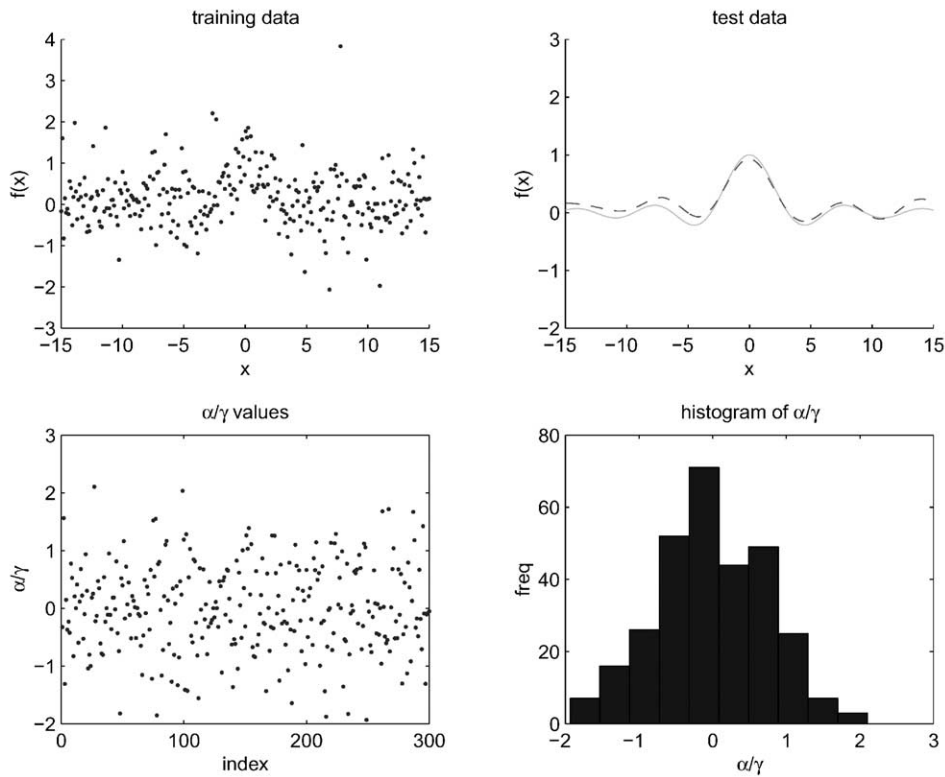


Fig. 5. (Continued) Weighted LS-SVM applied to the results from Fig. 4. The  $e_k$  distribution becomes Gaussian and the generalization performance on the test data improves.

a validation set or by means of cross-validation as illustrated in [33]) an update of  $(\gamma, \sigma)$  will be needed. The fact that this  $(\gamma, \sigma)$  determination can be kept ‘localized’ is a possible advantage of this method in comparison with other approaches which need to solve a QP problem for the several possible choices of the hyperparameters. Also note that the pruning method is in fact an iterative scheme where in each step one has to solve a KKT system. For interior point methods in solving standard quadratic programming type SVM’s [26], at each iteration step one solves a KKT system which has a similar form as in the LS-SVM case. However, in Algorithm 2 the size of the linear system that one solves decreases from step to step because  $N$  decreases. The convergence of Algorithm 2 for obtaining sparseness is guaranteed by construction.

## 5. Examples

In this example, we illustrate the method of weighted LS-SVM and sparse approximation. First, we show two examples of estimating a sinc function from noisy

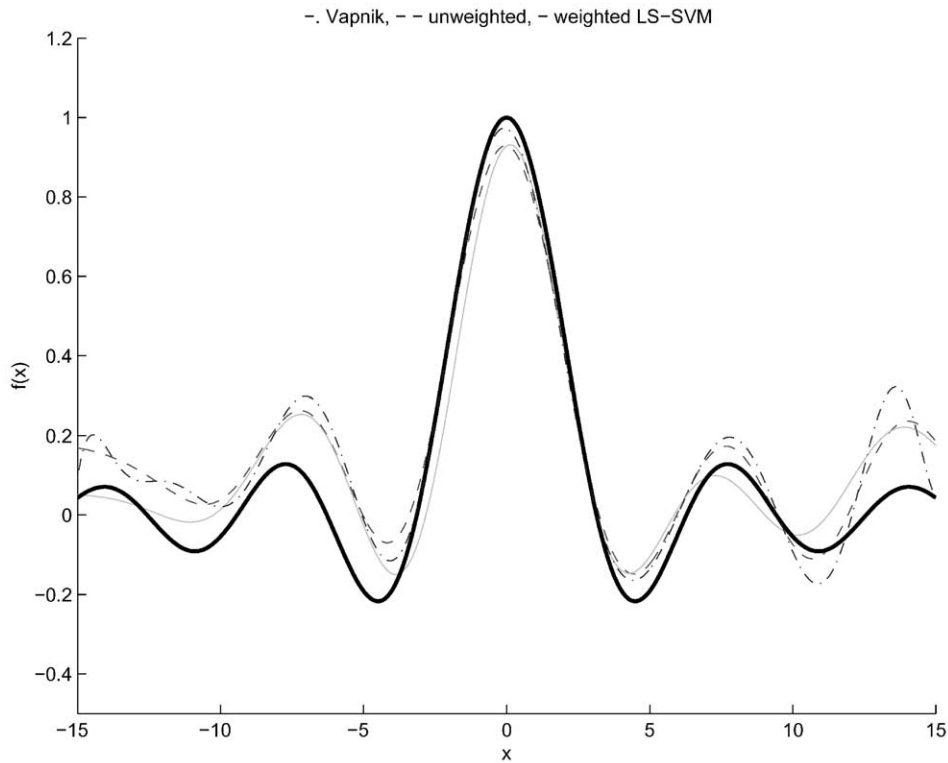


Fig. 6. (Continued) Comparison between standard SVM, unweighted LS-SVM and weighted LS-SVM. Weighted LS-SVM gives the best estimate for this example.

data: (a) strong outliers are superimposed on zero mean Gaussian noise distribution (Figs. 1 and 2); (b) non-Gaussian noise distribution (central  $t$ -distribution with 4 degrees of freedom, i.e. heavy tails [15]) (Figs. 4 and 5). Given is a training set of  $N = 300$  data points.

From the simulation results it is clear that in both cases the unweighted LS-SVM is quite robust and does not break down (Figs. 1 and 4). The generalization performance is further improved by applying weighted LS-SVM (Algorithm 1), shown in Figs. 2 and 5, respectively. The  $(\gamma, \sigma)$  hyperparameters are determined here by means of 10-fold cross-validation on the training data. The good generalization performance on fresh test data is shown for all cases. Figs. 2 and 5 show that the use of weighted LS-SVM leads to a  $e_k$  distribution which is closer to a Gaussian distribution, which leads to better estimates within the LS-SVM context.

An additional comparison with a standard SVM with Vapnik  $\varepsilon$ -insensitive loss function is made. The Matlab SVM Toolbox by Steve Gunn was used to generate the SVM results. Here  $\varepsilon = 0$  was taken and as upper bound on support values

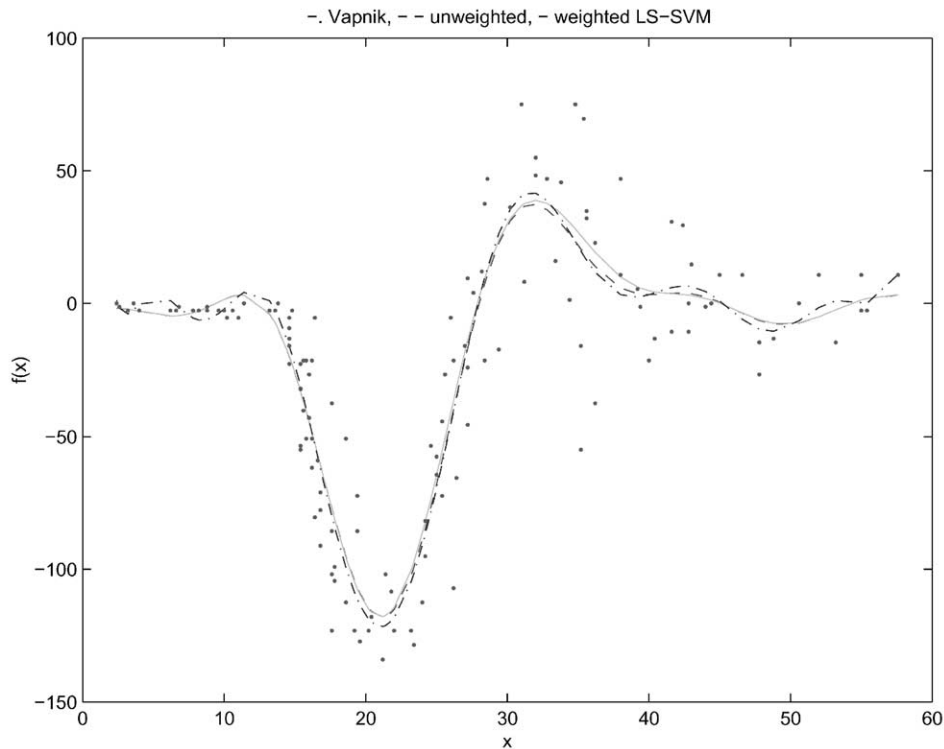


Fig. 7. Motorcycle dataset: comparison between standard SVM, unweighted LS-SVM and weighted LS-SVM.

$C = \text{Inf}$ . An optimal  $\sigma$  value was selected. Other  $\varepsilon, C$  combinations resulted in worse results. These comparative results are shown in Figs. 3 and 6. In these examples the weighted LS-SVM results show the best results. The unweighted LS-SVM is also quite robust. Due to the choice of a 2-norm this may sound surprising. However, one should be aware that only the output weights (support values  $\alpha$ ) follow from the solution to the linear system while  $(\gamma, \sigma)$  are to be determined at another level.

Figs. 7 and 8 show comparative results on the motorcycle data, a well-known benchmark data set in statistics [7]. The  $x$  values are time measurements in milliseconds after simulated impact and the  $y$  values are measurements of head acceleration. The  $x$  values are not equidistant and in some cases multiple  $y$  observations are present for certain  $x$  values. The data are heteroscedastic. In this sense it forms a challenging test case. Figs. 7 and 8 show the results from unweighted and weighted LS-SVM in comparison with standard SVM. In this example standard SVM suffers more from boundary effects. The tuning parameters are:  $\gamma = 2$ ,  $\sigma = 6.6$  (LS-SVM) and  $\sigma = 11$ ,  $\varepsilon = 0$ ,  $C = \text{Inf}$  (Vapnik SVM).

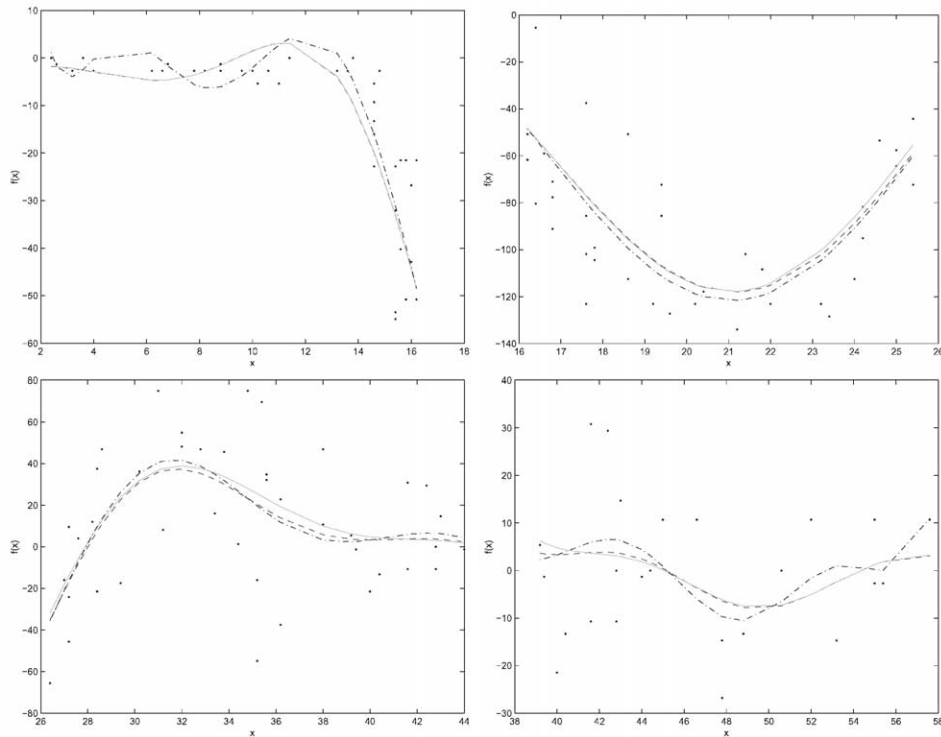


Fig. 8. (Continued) enlarged parts of the previous figure. From (Top-Left) it is visible that standard SVM suffers more from boundary effects with oscillatory behaviour. From (Bottom-Left) one observes that weighted LS-SVM improves in this region because of the correction for outliers.

Fig. 9 shows the improvements of weighted LS-SVM on the Boston housing data in comparison with unweighted LS-SVM. The weighted LS-SVM achieves an improved test set performance after determination of  $(\gamma, \sigma)$  by 10-fold CV on a randomly selected training set of 406 points. The remaining test set consisted of 100 points. The data were normalized except the binary variables. Optimal values of  $(\gamma, \sigma)$  were determined by 10-fold CV on the training set. The weighted LS-SVM resulted in a test set MSE error of 0.1638, which was an improvement over the unweighted LS-SVM test set MSE error of 0.1880. The improved performance is achieved by suppressing the outliers in the histogram shown in Fig. 9.

The sparse approximation procedure (Algorithm 2) is illustrated in Figs. 10 and 11. These results are obtained by starting from the training data of Figs. 1 and 2. The Figures show that robust estimates are obtained by a combination of weighted LS-SVM and sparse approximation by applying Algorithm 2. Some shiftings in the sorted support value spectrum are shown in Fig. 11. In this example about

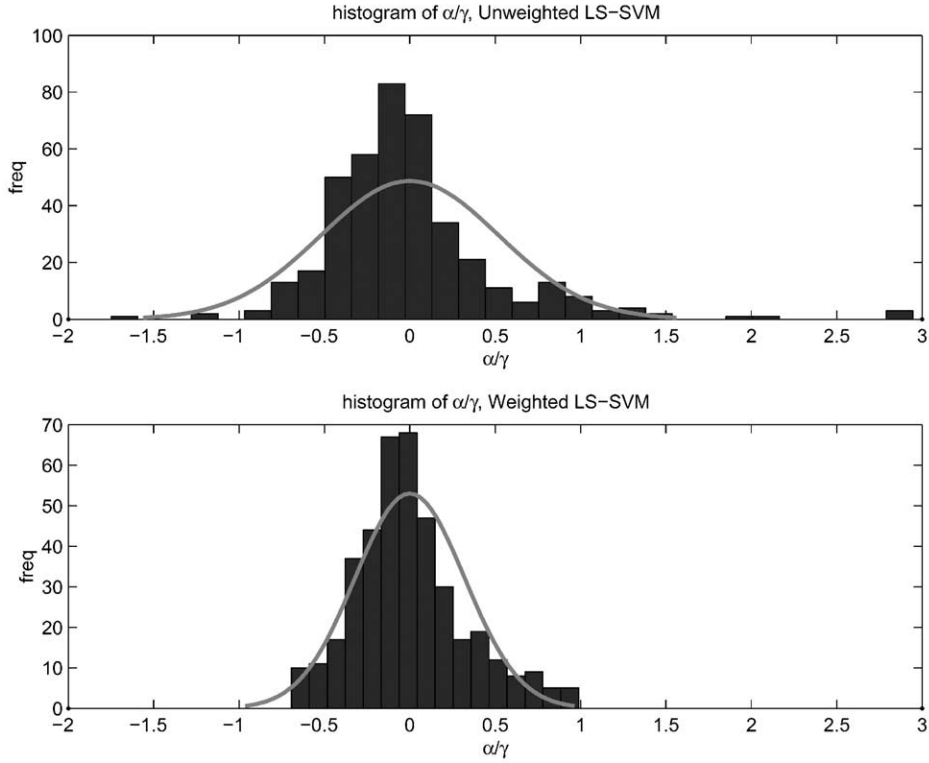


Fig. 9. Boston housing data set: (Top) histogram of  $e_k$  for unweighted LS-SVM with RBF kernel. The outliers are clearly visible; (Bottom) histogram resulting from weighted LS-SVM with improved test set performance.

20 support vectors are needed to guarantee a good generalization performance. As shown in [33] on many UCI data sets, optimizing the hyperparameters during the pruning process will improve the generalization performance. We want to stress here that pruning of the LS-SVM is in fact not needed, unless one wants to obtain a sparse representation of the original model.

## 6. Conclusions

We have shown how to obtain robust estimates within the LS-SVM framework in the case of outliers and heavy tailed non-Gaussian error distributions. This is done here by applying a weighted LS-SVM version. While least squares in standard parametric linear regression has a low breakdown point, LS-SVM's with RBF kernel have much better properties. Nevertheless, the robustness can be



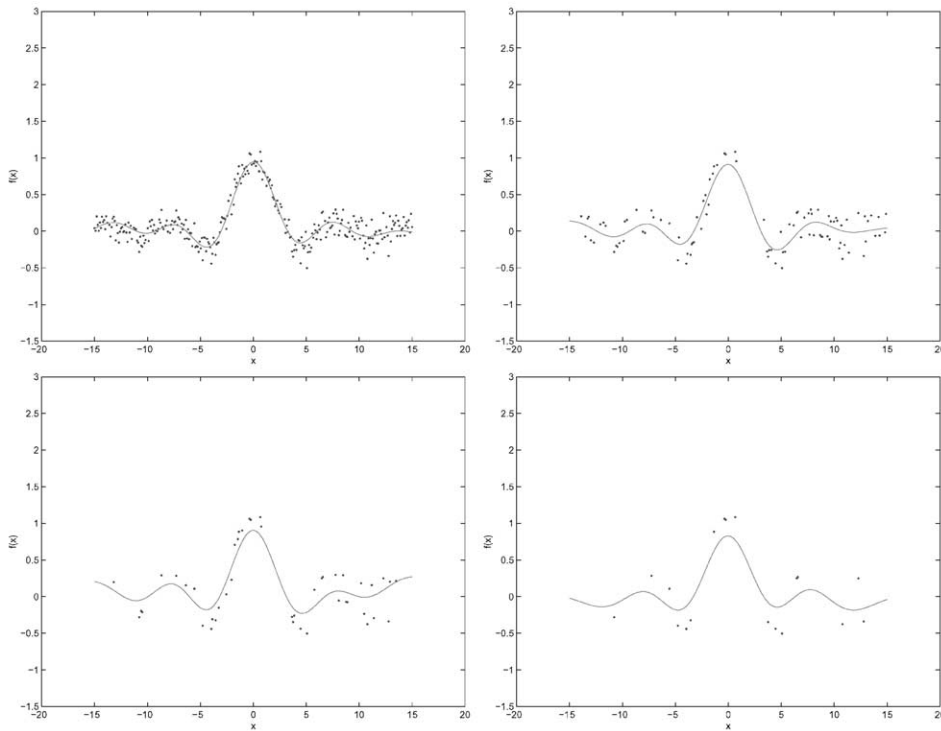


Fig. 10. Sparse approximation procedure of Algorithm 2, applied to the data of Figs. 1–2. Shown are some intermediate pruning results starting from 300 data points: 300 support vectors (SV) (Top-Left); 100 SV (Top-Right); 50 SV (Bottom-Left); 20 SV (Bottom-Right).

further enhanced by applying an additional weighted LS-SVM step. While standard SVM's possess a sparseness property for the solution vector, this feature is lost when considering LS-SVM's. However, we have shown how sparseness can be obtained by pruning the sorted support value spectrum if needed. The pruning procedure is based upon the solution vector itself which is an advantage in comparison with classical multilayer perceptron pruning. This procedure has the potential advantage of keeping the hyperparameter selection more localized. While standard SVM approaches start from choosing a given convex cost function and obtain a robust and sparse estimate in a top-down fashion, this procedure has the disadvantage that one should know in fact beforehand which cost function is statistically optimal. In this paper we have successfully demonstrated an alternative bottom-up procedure which starts from an un-weighted LS-SVM and then makes the solution robust by defining weightings based upon the error distribution. This provides motivation for further fundamental research in this direction for future work.

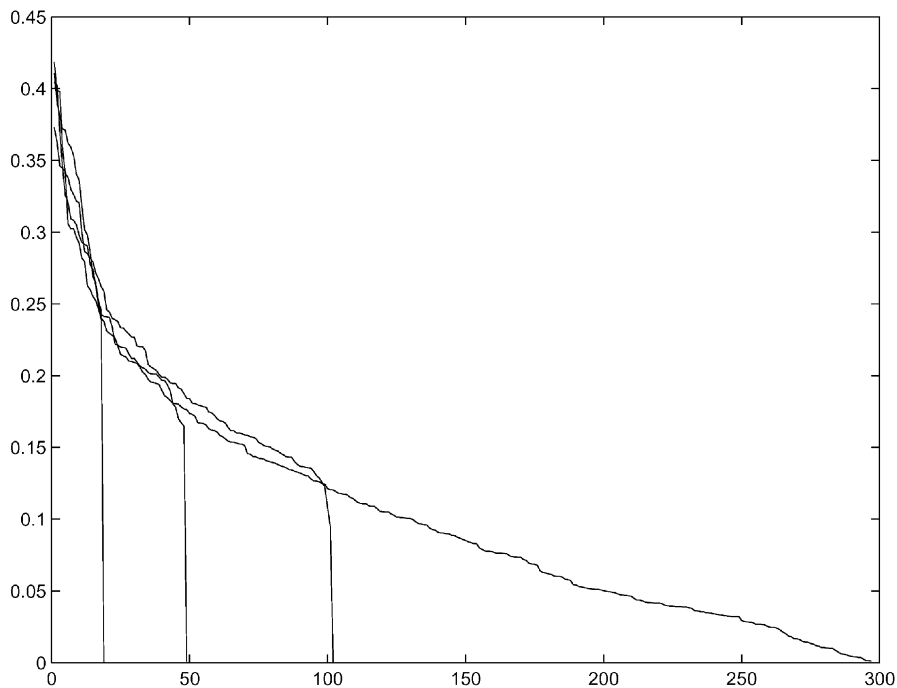


Fig. 11. Sorted support value spectrum related to the 4 steps shown in the previous Figure.

## References

- [1] D.F. Andrews, P.J. Bichel, F.R. Hampel, P.J. Huber, W.H. Rogers, J.W. Tukey, *Robust Estimates of Location: Survey and Advances*, Princeton University Press, Princeton, NJ, 1972.
- [2] C.M. Bishop, *Neural Networks for Pattern Recognition*, Oxford University Press, Oxford, 1995.
- [3] S.C. Chen, D.L. Donoho, M.A. Saunders, Atomic decomposition by basis pursuit, *SIAM J. Sci. Comput.* 20 (1) (1998) 33–61.
- [4] N. Cristianini, J. Shawe-Taylor, *An Introduction to Support Vector Machines*, Cambridge University Press, Cambridge, 2000.
- [5] H.A. David, Early sample measures of variability, *Statist. Sci.* 13 (4) (1998) 368–377.
- [6] D.L. Donoho, P.J. Huber, The notion of breakdown point, in: P. Bickel, K. Doksum, J.L. Hodges Jr. (Eds.), *A Festschrift for Erich Lehmann*, Wadsworth, Belmont, CA, 1983.
- [7] R.L. Eubank, *Nonparametric regression and spline smoothing*, *Statistics: textbooks and monographs*, Vol. 157, 2nd edition, Marcel Dekker, New York, 1999.
- [8] R. Fletcher, *Practical Methods of Optimization*, Wiley, New York, 1987.
- [9] F. Girosi, An equivalence between sparse approximation and support vector machines, *Neural Comput.* 10 (6) (1998) 1455–1480.
- [10] G.H. Golub, C.F. Van Loan, *Matrix Computations*, Johns Hopkins University Press, Baltimore, MD, 1989.
- [11] F.R. Hampel, E.M. Ronchetti, P.J. Rousseeuw, W.A. Stahel, *Robust Statistics: The Approach Based on Influence Functions*, Wiley, New York, 1986.

- [12] B. Hassibi, D.G. Stork, Second order derivatives for network pruning: optimal brain surgeon, in: S. Hanson, J. Cowan, L. Giles (Eds.), *Advances in Neural Information Processing Systems*, Vol. 5, Morgan Kaufmann, San Mateo, CA, 1993, pp. 164–171.
- [13] K. Hornik, M. Stinchcombe, H. White, Multilayer feedforward networks are universal approximators, *Neural Networks* 2 (1989) 359–366.
- [14] P.J. Huber, *Robust Statistics*, Wiley, New York, 1981.
- [15] N.L. Johnson, S. Kotz, *Distributions in Statistics: Continuous Univariate Distributions*, Vol. 1–2, Wiley, New York, 1970.
- [16] Y. Le Cun, J.S. Denker, S.A. Solla, Optimal brain damage, in: Touretzky (Ed.), *Advances in Neural Information Processing Systems*, Vol. 2, Morgan Kaufmann, San Mateo, CA, 1990, pp. 598–605.
- [17] T. Poggio, F. Girosi, Networks for approximation and learning, *Proc. IEEE* 78 (9) (1990) 1481–1497.
- [18] P.J. Rousseeuw, A. Leroy, *Robust Regression and Outlier Detection*, Wiley, New York, 1987.
- [19] P.J. Rousseeuw, B.C. van Zomeren, Unmasking multivariate outliers and leverage points, *J. Am. Statist. Assoc.* 85 (1990) 633–639.
- [20] C. Saunders, A. Gammerman, V. Vovk, Ridge regression learning algorithm in dual variables, *Proceedings of the 15th International Conference on Machine Learning (ICML'98)*, Morgan Kaufmann, 1998, pp. 515–521.
- [21] B. Schölkopf, K.-K. Sung, C. Burges, F. Girosi, P. Niyogi, T. Poggio, V. Vapnik, Comparing support vector machines with Gaussian kernels to radial basis function classifiers, *IEEE Trans. Signal Process.* 45 (11) (1997) 2758–2765.
- [22] B. Schölkopf, C. Burges, A. Smola (Eds.), *Advances in Kernel Methods—Support Vector Learning*, MIT Press, Cambridge, MA, 1998.
- [23] B. Schölkopf, S. Mika, C. Burges, P. Knirsch, K.-R. Müller, G. Rätsch, A. Smola, Input space vs. feature space in kernel-based methods, *IEEE Trans. Neural Networks* 10 (5) (1999) 1000–1017.
- [24] A. Smola, B. Schölkopf, K.-R. Müller, The connection between regularization operators and support vector kernels, *Neural Networks* 11 (1998) 637–649.
- [25] A. Smola, B. Schölkopf, On a kernel-based method for pattern recognition, regression, approximation and operator inversion, *Algorithmica* 22 (1998) 211–231.
- [26] A. Smola, *Learning with Kernels*, Ph.D. Thesis, GMD, Birlinghoven, 1999.
- [27] J.A.K. Suykens, J. Vandewalle (Eds.), *Nonlinear Modeling: Advanced Black-box Techniques*, Kluwer Academic Publishers, Boston, 1998.
- [28] J.A.K. Suykens, J. Vandewalle, Least squares support vector machine classifiers, *Neural Process. Lett.* 9 (3) (1999) 293–300.
- [29] J.A.K. Suykens, L. Lukas, P. Van Dooren, B. De Moor, J. Vandewalle, Least squares support vector machine classifiers: a large scale algorithm, *European Conference on Circuit Theory and Design, (ECCTD'99)*, PP. 839–842, Stresa Italy, August 1999.
- [30] J.A.K. Suykens, L. Lukas, J. Vandewalle, Sparse least squares support vector machine classifiers, *European Symposium on Artificial Neural Networks (ESANN 2000)*, Bruges Belgium, April 2000, pp. 37–42.
- [31] J.A.K. Suykens, J. Vandewalle, Recurrent least squares support vector machines, *IEEE Trans. Circuits Systems-I* 47 (7) (2000) 1109–1114.
- [32] J.A.K. Suykens, J. Vandewalle, B. De Moor, Optimal Control by Least Squares Support Vector Machines, *Neural Networks* 14 (1) (2001) 23–35.
- [33] T. Van Gestel, J.A.K. Suykens, B. Baesens, S. Viaene, J. Vanthienen, G. Dedene, B. De Moor, J. Vandewalle, Benchmarking least squares support vector machine classifiers, *Internal Report* 00-37, ESAT-SISTA, K.U.Leuven.
- [34] T. Van Gestel, J.A.K. Suykens, D. Baestaens, A. Lambrechts, G. Lanckriet, B. Vandaele, B. De Moor, J. Vandewalle, Financial time series prediction using least squares support vector machines within the evidence framework, *IEEE Trans. Neural Networks (special issue on Neural Networks in Financial Engineering)* 12 (4) (2001) 809–821.
- [35] V. Vapnik, *The Nature of Statistical Learning Theory*, Springer, New York, 1995.
- [36] V. Vapnik, *Statistical Learning Theory*, Wiley, New York, 1998.



**Johan A.K. Suykens** was born in Willebroek Belgium, May 18, 1966. He received the degree in Electro-Mechanical Engineering and the Ph.D. degree in Applied Sciences from the Katholieke Universiteit Leuven, in 1989 and 1995, respectively. In 1996 he was a Visiting Postdoctoral Researcher at the University of California, Berkeley. At present, he is a Postdoctoral Researcher with the Fund for Scientific Research FWO Flanders. His research interests are mainly in the areas of the theory and application of nonlinear systems and neural networks. He is author of the book “Artificial Neural Networks for Modelling and Control of Non-linear Systems” and editor of the book “Nonlinear Modeling: Advanced Black-Box Techniques”, published by Kluwer Academic Publishers. The latter resulted from an International Workshop on Nonlinear

Modelling with Time-series Prediction Competition that he organized in 1998. He has served as associate editor for the IEEE Transactions on Circuits and Systems-I (1997–1999) and since 1998 he is serving as associate editor for the IEEE Transactions on Neural Networks. He received a Best Paper Award as Winner for Best Theoretical Developments in Computational Intelligence at ANNIE’99 and an IEEE Signal Processing Society 1999 Best Paper Award, for his contributions on NLq Theory. He is a recipient of the International Neural Networks Society INNS 2000 Young Investigator Award for significant contributions in the field of neural networks.



**Jos De Brabanter** received the diploma of Ind. Ing. (Brussels, Belgium) in 1990, the Master degree in Artificial Intelligence in 1996 and the Master degree of Statistics in 1997, both from the K.U. Leuven (Belgium). He is currently pursuing the Doctoral degree in Applied Sciences at the Department of Electrical Engineering of the K.U. Leuven. His scientific research interests are in the area of statistics and its application to neural networks and support vector learning.



**Lukas** received the B.Eng. degree in Computer Engineering in 1995 from the Institute of Technology, Bandung (ITB), Indonesia, and the Master degree in Artificial Intelligence from the K.U. Leuven, Belgium in 1998. He is currently pursuing the Doctoral degree at the Department of Electrical Engineering at the K.U. Leuven. His research interests are in the area of support vector machines, biomedical engineering and artificial intelligence.



**Joos Vandewalle** was born in Kortrijk, Belgium, in August 1948. He obtained the electrical engineering degree and doctoral degree in applied sciences, both from the Katholieke Universiteit Leuven, Belgium in 1971 and 1976, respectively. From 1976 to 1978 he was Research Associate and from July 1978 to July 1979, he was Visiting Assistant Professor both at the University of California, Berkeley. Since July 1979 he is back at the ESAT Laboratory of the Katholieke Universiteit Leuven, Belgium, where he is Full Professor since 1986. He is an Academic Consultant since 1984 at the VSDM group of IMEC (Interuniversity Microelectronics Center, Leuven). From August 1996 to August 1999 he was Chairman of the Department of Electrical Engineering at the Katholieke Universiteit Leuven. Since August 1999 he is the vice-dean of the Department of Engineering. He teaches courses in linear algebra, linear

and nonlinear systems and circuit theory, signal processing and neural networks. His research interests are mainly in mathematical system theory and its applications in circuit theory, control, signal processing, cryptography and neural networks. He has authored and co-authored more than 200 papers in these areas. He is co-author with S. Van Huffel of the book “The Total Least Squares Problem” and co-editor with T. Roska of the book “Cellular Neural Networks”. He is a member of the editorial board of “Journal A”, a Quarterly Journal of Automatic Control and of the International Journal of Circuit Theory and its Applications, Neurocomputing and the Journal of Circuit Systems and Computers. From 1989 till 1991, he was associate editor at the IEEE Transactions on Circuits and Systems in the area of nonlinear and neural networks. He was elected fellow of IEEE in 1992 for contributions to nonlinear circuits and systems. During 1991–1992 he held the Francqui chair of Artificial Neural Networks at the University of Liège, Belgium. He is one of the three coordinators of the Interdisciplinary Center for Neural Networks ICNN that was set up in 1993 in order to stimulate the interaction and cooperation among the 50 researchers on neural networks at the Katholieke Universiteit Leuven.