

# A Robust LS-SVM Regression

József Valyon, and Gábor Horváth

**Abstract**—In comparison to the original SVM, which involves a quadratic programming task; LS-SVM simplifies the required computation, but unfortunately the sparseness of standard SVM is lost. Another problem is that LS-SVM is only optimal if the training samples are corrupted by Gaussian noise. In Least Squares SVM (LS-SVM), the nonlinear solution is obtained, by first mapping the input vector to a high dimensional kernel space in a nonlinear fashion, where the solution is calculated from a linear equation set. In this paper a geometric view of the kernel space is introduced, which enables us to develop a new formulation to achieve a sparse and robust estimate.

**Keywords**—Support Vector Machines, Least Squares Support Vector Machines, Regression, Sparse approximation.

## I. INTRODUCTION

THIS paper focuses on the least squares version of SVM [1], the LS-SVM [2], whose main advantage is that it is computationally more efficient than the standard SVM method. In this case training requires the solution of a linear equation set instead of the long and computationally hard quadratic programming problem involved by the standard SVM.

The method effectively reduces the algorithmic complexity, however for really large problems, comprising a very large number of training samples, even this least-squares solution can become highly memory and time consuming.

Whereas the least squares version incorporates all training data in the network to produce the result, the traditional SVM selects some of them (the support vectors) that are important in the regression. The sparseness of traditional SVM can also be reached with LS-SVM by applying a pruning method [3][4]. Unfortunately if the traditional LS-SVM pruning method is applied, the performance declines proportionally to the eliminated training samples, since the information (input-output relation) they described is lost. Another problem is that this iterative method multiplies the algorithmic complexity.

The training data is often corrupted by noise, which –if not handled properly– misleads the training. Another modification of the method, called weighted LS-SVM [1][5], is aimed at reducing the effects of non-Gaussian noise (e.g. outliers).

The biggest problem is that pruning and weighting –

This work was partly sponsored by National Fund for Scientific Research (OTKA) under contract T 046771.

J. Valyon is with Budapest University of Technology and Economics-Department of Measurement and Information Systems, Budapest, Hungary, H-1521, pf. 91. (phone +36 1 463-2057; fax +36 1 463-4112; e-mail: valyon@mit.bme.hu).

G. Horváth (e-mail: horvath@mit.bme.hu).

although their goals do not rule out each other– cannot be used at the same time, because they work in opposition. The generalized approach presented in this paper enables us to accomplish both goals by allowing a more universal construction and solution of the LS-SVM equation set.

This paper proposes a geometric view of the kernel space and the linear solution that is based on the mapped training samples.

In the LS-SVM solution, the training samples are mapped to a kernel space, where a hyperplane is fitted on these points. In this case all the training samples are used to achieve a result, which consequently isn't sparse. To trade off between training error and a smooth solution a regularization parameter is used, which is the same for all samples, and can be considered as a predefined, intentional error term in the kernel space fitting.

Our proposition is to use a kernel space of smaller dimensionality, which means that by mapping all training samples, the hyperplane can be fitted many ways, since there are more equations than unknown. This means that the position of the hyperplane and consequently the training errors –the distances from this plane– can be automatically determined according to the distribution of many mapped points.

The LS-SVM method is capable of solving both classification and regression problems. The classification approach is easier to understand and more historic. The present study concerns regression, but it must be emphasized that all presented methods can be applied to classification as well.

This paper is organized as follows. Before going into the details LS-SVM, and its extensions of *pruning* and *weighting* are summarized in Section II. Section III provides a geometric interpretation of the kernel space and summarizes the main idea behind the propositions. Section IV contains the details of the solution: It shows how partial reduction is used to achieve an overdetermined equation set, and proposes robust solutions for this.

Section V. contains some experimental results, while in section VI. the conclusions are drawn.

## II. A BRIEF INTRODUCTION TO THE LS-SVM METHOD

Given the  $\{\mathbf{x}_i, d_i\}_{i=1}^N$  training data set, where  $\mathbf{x}_i \in \mathcal{R}^p$  represents a  $p$ -dimensional input vector and  $d_i = y_i + z_i$ ,  $d_i \in \mathcal{R}$  is a scalar measured output, which represents the  $y_i$  system output corrupted by some  $z_i$  noise. Our goal is to

construct an  $y = f(\mathbf{x})$  function, which represents the dependence of the output  $y_i$  on the input  $\mathbf{x}_i$ . Let's define the form of this function as formulated below:

$$y = \sum_{i=1}^h w_i \varphi_i(\mathbf{x}) + b = \mathbf{w}^T \boldsymbol{\varphi}(\mathbf{x}) + b, \quad (1)$$

$$\mathbf{w} = [w_1, w_2, \dots, w_h]^T, \quad \boldsymbol{\varphi} = [\varphi_1, \varphi_2, \dots, \varphi_h]^T.$$

The  $\boldsymbol{\varphi}(\cdot) : \mathcal{R}^p \rightarrow \mathcal{R}^h$  is a mostly non-linear function, which maps the data into a higher –possibly infinite– dimensional feature space. The main difference from the standard SVM is that LS-SVM involves equality constraints instead of inequality ones and works with a least squares cost function [1]. The optimization problem and the inequality constraints are defined by the following equations ( $i = 1, \dots, N$ ):

$$\min_{\mathbf{w}, b, e} J_p(\mathbf{w}, e) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \frac{1}{2} \sum_{i=1}^N e_i^2 \quad (2)$$

$$\text{with constraints: } d_i = \mathbf{w}^T \boldsymbol{\varphi}(\mathbf{x}_i) + b + e_i.$$

The  $C \in \mathcal{R}^+$  is the trade-off parameter between a smoother solution, and training errors. From this, a Lagrangian is formed

$$L(w, b, e; \alpha) = J_p(w, e) - \sum_{i=1}^N \alpha_k \{ \mathbf{w}^T \boldsymbol{\varphi}(\mathbf{x}_i) + b + e_i - d_i \}. \quad (3)$$

The solution concludes in a constrained optimization, where the conditions for optimality lead to the following overall solution:

$$\begin{bmatrix} 0 & \bar{\mathbf{1}}^T \\ \bar{\mathbf{1}} & \boldsymbol{\Omega} + C^{-1} \mathbf{I} \end{bmatrix} \begin{bmatrix} b \\ \boldsymbol{\alpha} \end{bmatrix} = \begin{bmatrix} 0 \\ \mathbf{d} \end{bmatrix}, \quad (4)$$

$$d = [d_1, d_2, \dots, d_N], \quad \boldsymbol{\alpha} = [\alpha_1, \alpha_2, \dots, \alpha_N],$$

$$\bar{\mathbf{1}} = [1, \dots, 1], \quad \Omega_{i,j} = K(\mathbf{x}_i, \mathbf{x}_j) = \boldsymbol{\varphi}^T(\mathbf{x}_i) \boldsymbol{\varphi}(\mathbf{x}_j),$$

where  $K(\mathbf{x}_i, \mathbf{x}_j)$  is the kernel function, and  $\boldsymbol{\Omega}$  is the kernel matrix. The result is:

$$y = \sum_{i=1}^N \alpha_i K(\mathbf{x}, \mathbf{x}_i) + b \quad (6)$$

A detailed description of LS-SVM can be found in refs. [2]-[5].

**LS-SVM pruning** - One of the main drawbacks of the least-squares solution is that, it is not sparse, because –unlike the original SVM– it incorporates all training vectors in the result. In order to get a sparse solution, a pruning method must be used. Since the  $|\alpha_i|$  support values are proportional to the errors at data points:

$$\alpha_i = C e_i, \quad (7)$$

the irrelevant points are left out, by iteratively leaving out the least significant vectors. These are the ones corresponding to the smallest  $|\alpha_i|$  values.

In the case of the classical SVM sparseness is achieved by the use of such loss functions, where errors smaller than  $\varepsilon$  are ignored (e.g.  $\varepsilon$ -insensitive loss function). This method reduces the difference between SVM and LS-SVM, because the omission of some data points implicitly corresponds to creating an  $\varepsilon$ -insensitive zone [1].

The described method leads to a sparse model, but some questions arise: How many neurons are needed in the final model? How many iterations it should take to reach the final model?

Another problem is that a usually large linear system must be solved in each iteration. Pruning is especially important if the number of training vectors is large. In this case however, the iterative method is not very effective.

**Weighted LS-SVM** - This method addresses the problem of noisy data –like outliers in a dataset–, by using a weighting factor in the calculation based on the error variables determined from a previous –first an unweighted– solution. The method uses a bottom-up approach by starting from a standard solution, and calculating one or more weighted LS-SVM based on the previous result. The weighted LS-SVM is formed as:

$$\begin{bmatrix} 0 & \bar{\mathbf{1}}^T \\ \bar{\mathbf{1}} & \boldsymbol{\Omega} + \mathbf{V}_C \end{bmatrix} \begin{bmatrix} b \\ \boldsymbol{\alpha} \end{bmatrix} = \begin{bmatrix} 0 \\ \mathbf{d} \end{bmatrix} \quad (8)$$

where

$$\mathbf{V}_C = \text{diag} \left\{ \frac{1}{C v_1}, \dots, \frac{1}{C v_N} \right\}. \quad (9)$$

The  $v_i$  weighting is designed such, that the results improve in view of robust statistics. Large  $e_i$ -s mean a small weight and vice versa.

A common property of the described methods is that they are all iterative, where every step is based on the result of an LS-SVM learning. This means, that the entire large problem must be solved at least once, and a relatively large one in every further iteration step. Another drawback is that pruning and weighting cannot be easily combined, because the methods favor contradictory types of points. While pruning drops the training points belonging to small  $\alpha_i$ -s, the weighted LS-SVM increases the effects of these points.

### III. THE MAIN IDEA

When an LS-SVM is constructed from  $N$  training samples:

1. The samples are mapped to an  $N+1$  dimensional kernel space, where  $N$  dimensions are defined by the kernel functions and one is the desired output.
2. A hyperplane is fitted on these mapped samples. The hyperplane is determined by  $N$  mapped points, and one additional constraint ( $\sum_{i=1}^N \alpha_i = 0$ ).

The approximated answer for a new sample results from mapping it into  $N$  dimensions, and calculating the corresponding point on the fitted hyperplane (dimension  $N+1$ ). Therefore in case of this solution, we expect that when a new sample is transformed to this kernel space, the desired output will be close to this hyperplane. For the sake of generalization and to avoid overfitting the accuracy of the fit can be adjusted through regularization (the  $C$  hyperparameter).

Some questions still need answers: How many –and which–

dimensions are needed in the kernel space? How good is my approximation, for a specific mapping? What is a good value of  $C$ , or more generally, how should the hyperplane be placed in the kernel space?

Our proposition is to take control of the problem in the kernel space by:

- ◆ controlling the dimensionality of this space (see IV. A.),
- ◆ finding a better linear hyperplane in the kernel space (see IV. B.),
- ◆ choosing an appropriate kernel space (IV. C.).

Having fewer dimensions in the kernel space, results in a sparse solution, while at the same time it increases the number of mapped points that can be used to determine the linear fit. Having more points than dimensions in the kernel space allows us to optimize the linear fit. The dimensionality of the kernel space high is enough, if samples (not used in determining this space) fall close to this plane after mapping (see Fig. 1.).

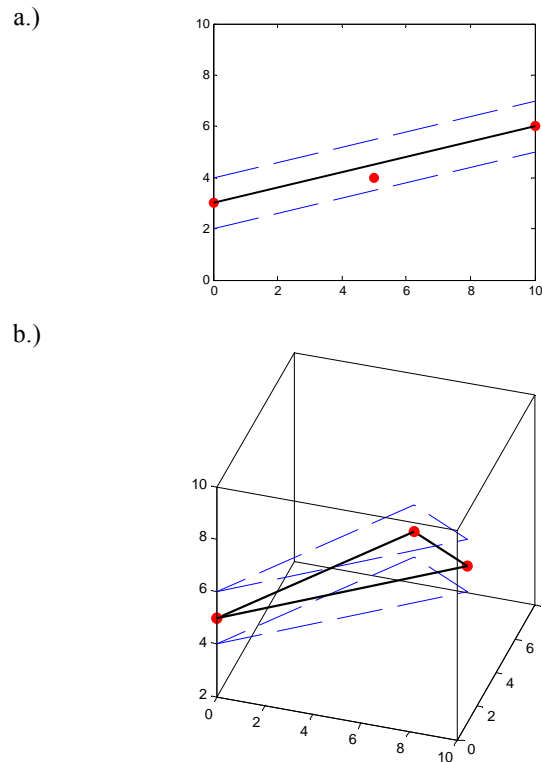


Fig. 1 The image of training samples in a kernel space of different dimensions. Using all three samples as support vectors (kernel centers), a three dimensional kernel can space guarantee exact fit for the samples. The dashed lines represent a zone in which errors can be accepted (corresponding to the  $\epsilon$ -insensitivity of SVM)

#### IV. THE PROPOSED METHODS

This section proposes some modifications and extensions to the standard LS-SVM. Their purpose is to gain control over network size, to reduce complexity and to improve the quality of the results.

##### A. Using an Overdetermined Equation Set

If the training set consists of  $N$  samples, then our original linear equation set (see eq. 5) will have  $(N+1)$  unknowns, the  $\alpha_i$ -s and  $b$ ,  $(N+1)$  equations and  $(N+1)^2$  multipliers. These factors are mainly the values of the  $K(\mathbf{x}_i, \mathbf{x}_j)$  kernel function calculated for every combination of the training input pairs. The cardinality of the training set therefore determines the size of the kernel matrix, which plays a major part in the solution, as algorithmic complexity; the complexity of the result etc. depends on this.

To reduce the equation set, columns and/or rows may be omitted.

- ◆ If the  $k$ -th column is left out, then the corresponding  $\alpha_k$  is also deleted, therefore the resulting model will be smaller. The  $\sum_{i=1}^N \alpha_i = 0$  condition automatically adapts, since the remaining  $\alpha$ -s will still add up to zero.
- ◆ If the  $j$ -th row is deleted, then the condition defined by the  $(\mathbf{x}_j, d_j)$  training sample is lost, because the  $j$ -th equation is removed.

The most important component of the main matrix is the  $\Omega$  kernel matrix; its elements are the results of the kernel function for pairs of training inputs:

$$\Omega_{i,j} = K(\mathbf{x}_i, \mathbf{x}_j) \quad (10)$$

To reduce the size of  $\Omega$  some training samples should be omitted. Each column of the kernel matrix represents an additive term in the final solution, with a kernel function centered on the corresponding  $\mathbf{x}_i$  input. The rows however, represent the input-output relations, described by the training points. It can be seen that in order to reach sparseness the number of columns must be reduced. The following reduction techniques can be used on the kernel matrix (the names of these techniques are introduced here for easier discussion):

##### Traditional full reduction

A training sample  $(\mathbf{x}_k, d_k)$  is fully omitted, therefore both the column and the row corresponding to this sample are eliminated. In this case however reduction also means that the knowledge represented by the numerous other samples are lost. This is exactly the case in traditional LS-SVM pruning since pruning iteratively omits some training points. The information embodied in these points is entirely lost.

To avoid this information loss, one may use the technique referred here as partial reduction.

##### The proposed partial reduction

In partial reduction, the omission of a training sample  $(\mathbf{x}_k, d_k)$  means that only the corresponding column is eliminated, while the row –which defines an input-output relation– is kept. Eliminating the  $k$ -th column reduces the model complexity, while keeping the  $k$ -th row means that the weighted sum of that row should still meet the  $d_k$  regression goal (as closely as possible).

By selecting some (e.g.  $M$ ,  $M < N$ ) vectors as “support



so the weights can be calculated from the statistical properties of the points (distribution of many points) in the kernel space. Another important difference is that the proposed weighted solution is also *sparse*.

- ◆ CUSTOM WEIGHTING - this method can be used if one had *a priori* knowledge about the quality of the samples. If so, weights can be defined to determine how much each learning sample influences the fit. Samples known to have less noise are expected to fit more, than low-quality ones.

The weights should transform the response variances to a constant value. If the variances of the data are known, the weights are given by:

$$w_i = 1/\sigma_i^2 . \quad (19)$$

- ◆ BISQUARE WEIGHTS – a method that minimizes a weighted sum of squares, where the weight of each data point depends on its distance from the fitted line. The farther away is the point, the less weight it gets. This method fits the hyperplane to the bulk of the data with the least squares approach, while it minimizes the effect of outliers (Fig. 2.). More details on robust regression can be found in [9][10].

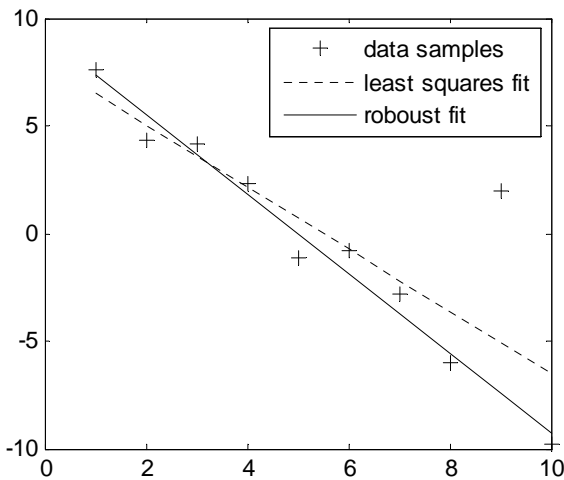


Fig. 2 The least squares and the robust (bisquare) fitting in two dimensions

### C. Selecting Support Vectors

Standard SVM automatically selects the support vectors. To achieve sparseness by partial reduction, the linear equation set has to be reduced in such a way, that the solution of this reduced (overdetermined) problem is the closest to what the original solution would be.

As the matrix is formed from columns we can select a linearly independent subset of column vectors and omit all others, which can be formed as linear combinations of the selected ones. This can be done by finding a “basis” (the quote indicates, that this basis is only true under certain conditions defined later) of the coefficient matrix, because the basis is by definition the smallest set of vectors that can solve the problem. The linear dependence discussed here, does not

mean exact linear dependence, because the method uses an adjustable tolerance value when determining the “resemblance” (parallelism) of the column vectors. The use of this tolerance value is essential, because none of the columns of the coefficient matrix will likely be exactly dependent (parallel). The reduction is achieved as a part of transforming the  $A^T$  matrix into reduced row echelon form, using a slight modification of Gauss-Jordan elimination with partial pivoting [11]. This method returns a list of the column vectors which are linearly independent from the others considering a tolerance  $\varepsilon'$ .

The tolerance ( $\varepsilon'$ ) can be related to the  $\varepsilon$  parameter of the standard SVM, because it has similar effects. The larger the tolerance, the fewer vectors the algorithm will select. If the tolerance is chosen too small, than a lot of vectors will seem to be independent, resulting in a larger network. As stated earlier the standard SVM’s sparseness is due to the  $\varepsilon$ -insensitive zone, which allows the samples falling inside this boundary to be neglected. According to this, it may not be very surprising to find that an additional parameter is needed to achieve sparseness in LS-SVM, and this parameter corresponds to the one, which was originally left when changing from the SVM to the standard least squares solution.

The basic idea of doing a feature selection in the kernel space is not new. The nonlinear principal component analysis technique, the Kernel PCA uses a similar idea [12]. A basis selection from the kernel matrix has been shown in [13].

## V. EXPERIMENTS

The next figures show the results for a simple illustrative experiment, the *sinc(x)* regression. The training set contains 50 data samples corrupted with Gaussian noise.

Fig. 3 shows the results of custom weighting. We have 60 samples with additive Gaussian noise, where the  $\sigma$  of the noise is known for all samples. It can be seen, that the effect of noise is reduced. The original LS-SVM is plotted, because the Weighted LS-SVM would give almost the same results as the partially reduced solution, but in this case we have a sparse solution.

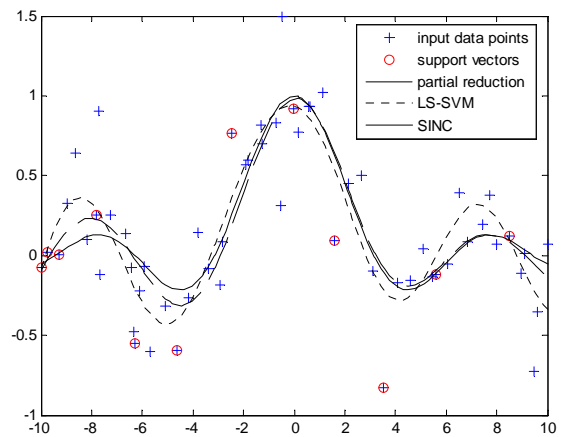


Fig. 3 Custom weighting is applied with partial reduction (The

LS-SVM is not weighted.)

The following experiment (Fig. 4) shows the same problem as Fig. 3, but in this case a few data points are corrupted to provide outliers.

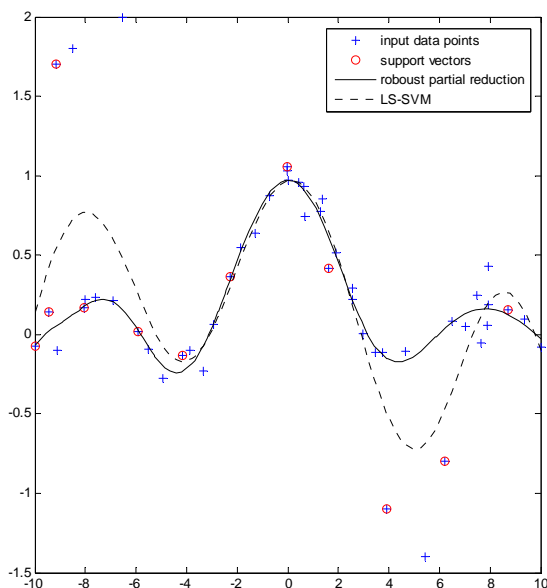


Fig. 4 The continuous black line plots the result for a partially reduced LS-SVM solved by the bisquare weights method. The dashed line is the original LS-SVM

It can be seen that by using a robust bisquare fitting, the effect of the outliers was successfully reduced.

It is important to mention that the result of the LS<sup>2</sup>-SVM is sparse, consisting of only 12 support vectors. If the number of training samples is very high for the problem complexity, than the gain in the network size can be rather large.

## VI. CONCLUSION

In this paper a geometric view and a generalized formulation of the least squares support vector machine was presented. The basic idea is that by reducing the dimensionality of the kernel space, the hyperplane fitted to the mapped training samples can be optimized according to their distribution. This is especially important, to deal with non Gaussian noise.

The described solution achieves two important results simultaneously:

- ◆ a sparse LS-SVM solution,
- ◆ the effect of noise is reduced.

## REFERENCES

[1] V. Vapnik, "The Nature of Statistical Learning Theory", New-York: Springer-Verlag, 1995  
 [2] J. A. K. Suykens, V. T. Gestel, J. De Brabanter, B. De Moor, J. Vandewalle, "Least Squares Support Vector Machines", World Scientific, 2002

[3] J. A. K. Suykens, L. Lukas, and J. Vandewalle, "Sparse approximation using least squares support vector machines", *IEEE International Symposium on Circuits and Systems ISCAS'2000*, 2000  
 [4] J. A. K. Suykens, L. Lukas, and J. Vandewalle, "Sparse least squares support vector machine classifiers", *ESANN'2000 European Symposium on Artificial Neural Networks*, 2000, pp. 37-42.  
 [5] J.A.K. Suykens, J. De Brabanter, L. Lukas, and J. Vandewalle, "Weighted least squares support vector machines: robustness and sparse approximation", *Neurocomputing*, 2002, pp. 85-105  
 [6] Yuh-Jye Lee and O. L. Mangasarian, "RSVM: Reduced Support Vector Machines", *Proceedings of the First SIAM International Conference on Data Mining*, Chicago, 2001, April 5-7.  
 [7] J. Vallyon and G. Horváth, "A generalized LS-SVM", *SYSID'2003 Rotterdam*, 2003, pp. 827-832.  
 [8] J. Vallyon and G. Horváth, "A Sparse Least Squares Support Vector Machine Classifier", *Proceedings of the International Joint Conference on Neural Networks IJCNN 2004*, 2004, pp. 543-548.  
 [9] Holland, P. W., and R. E. Welsch, "Robust Regression Using Iteratively Reweighted Least-Squares," *Communications in Statistics: Theory and Methods*, A6, 1977, pp. 813-827.  
 [10] Huber, P. J., *Robust Statistics*, Wiley, 1981.  
 [11] W. H. Press, S. A. Teukolsky, W. T. Wetterling and B. P. Flannery, "Numerical Recipes in C", Cambridge University Press, Books On-Line, Available: www.nr.com, 2002  
 [12] B. Schölkopf, S. Mika, C.J.C. Burges, P. Knirsch, K.-R. Müller, G. Rätsch, and A. Smola, "Input space vs. feature space in kernel-based methods". *IEEE Transactions on Neural Networks*, 1999, 10(5), pp. 1000-1017.  
 [13] G. Baudat and F. Anouar, "Kernel-based methods and function approximation". In *International Joint Conference on Neural Networks*, pages 1244-1249, Washington DC, 2001, July 15-19.  
 [14] H. Golub and Charles F. Van Loan, *Matrix Computations*, Gene Johns Hopkins University Press, 1989.