## F-SVR: A NEW LEARNING ALGORITHM FOR SUPPORT VECTOR REGRESSION

*Mireille Tohmé*<sup>(1)</sup>, *Régis Lengellé*<sup>(2)</sup>

<sup>(1)</sup> FORENAP Frp, 27, Rue du 4eme RSM , 68250 Rouffach, France
<sup>(2)</sup> ICD-LM2S (FRE CNRS 2848), Troyes University of Technology, BP 2060, 10010 Troyes cedex, France

## ABSTRACT

In this paper, we present a new method for optimizing Support Vector Machines for regression problems. This algorithm searches for efficient feasible directions. Within these selected directions, we choose the best one, i.e. the one, coupled with an optimal step analytical evaluation, that ensures a maximum increase of the objective function. The resulting solution, the gradient and the objective function are recursively determined and the Gram matrix has not to be stored. Our algorithm is based on SVM-Torch proposed by Collobert for regression, which is similar to SVM-Light suggested by Joachims for classifications problems, but adapted to regression problems. We are also inspired by LASVM proposed by Bordes for classification problems. F-SVR algorithm uses a new efficient working set selection heuristic, ingeniously exploits quadratic function properties, so it is fast as well as easy to implement and is able to perform on large data sets.

*Index Terms*— Training, algorithm, Support Vector Machines

## 1. INTRODUCTION

Over the past years, Support Vector Machine (SVM) have become an important tool for solving pattern-recognition and regression problems. Recently, Smola and Scholkopf [1] proposed an iterative algorithm, called Sequential Minimal Optimization (SMO), for solving the regression problem using SVM. Their decomposition algorithm used an analytical solution for the subproblems, and they proposed to select two pairs of variables as a working set. This algorithm is an extension of the SMO proposed by Platt [2] for SVM classification. Bordes [3] suggested LASVM which is a reorganization of the SMO sequential direction search method and, as such, converges to the solution of the SVM quadratic programming problem. Collobert et al [4] used a decomposition algorithm similar to the one proposed by Joachims [5]. They suggested to select two variables instead of two pairs of variables as did Scholkopf and Smola. In fact, working with pairs of variables would force the algorithm to do many computations with null variables until the end of the optimization process. In this paper, we take advantage of the decomposition method to present a new learning algorithm for solving regression problems. F-SVR relies on the following four points:

- 1. The particular analytical expression of the equality constraint of the optimization problem
- 2. An efficient and new heuristic for selecting a possible working set
- 3. The analytical determination of the optimal step-size of the gradient algorithm which, coupled with the previous heuristic, is new within the context of SVM training
- 4. The analytical recursive computation of the solution, the gradient and the objective function

The two last properties result from the quadratic nature of the objective function, ensure the efficacy of F-SVR algorithm and make the difference with traditional algorithms. We begin this paper with a brief review of the basic technique used for implementing Support Vector Machines for regression problems. Then we give an overview of F-SVR by introducing in details its different steps. Section 4 presents the experimental results and we finally conclude with some perspectives.

#### 2. SUPPORT VECTOR REGRESSION

In this section, we concisely review the fundamentals of Support Vector Regression in the non linear case. As usual, suppose we are given training data  $(x_1, y_1), \dots, (x_N, y_N) \subset \mathcal{X} \times \mathbb{R}$  where  $\mathcal{X}$  denotes the space of the input patterns, for example  $\mathbb{R}^d$ . In  $\varepsilon$ -SV regression [6], Vapnik introduced the  $\varepsilon$ -insensitive loss function that enables sparsity within the Support Vectors (SV). Our goal is to find a function f(x) that has at most  $\varepsilon$  deviation from the actually obtained targets  $y_i$  for all the training data, and at the same time is as smooth as possible. In other words, we do not care about errors as long as they are less than  $\varepsilon$ , but will pay for any deviation larger than  $\varepsilon$ . So we need to construct a regression function

$$f(x) = \langle w, \Phi(x) \rangle + b \qquad w \in \Phi(\mathcal{X}), b \in \mathbb{R}$$
(1)

on a feature space  $F = \Phi(\mathcal{X})$ . Here, w is a vector in F, and  $\Phi(x)$  maps the input x to a vector in F. The w and b in (1) are obtained by solving the following quadratic optimization

problem:

$$\begin{array}{l} \underset{w,b,\xi,\xi^{*}}{\operatorname{minimize}} \quad \frac{1}{2} \|w\|^{2} + C \sum_{i=1}^{N} \left(\xi_{i} + \xi_{i}^{*}\right) \\ \text{subject to} \begin{cases} y_{i} - \left(\left\langle w, \Phi\left(\mathbf{x}\right)\right\rangle + b\right) \leq \varepsilon + \xi_{i} \\ \left(\left\langle w, \Phi\left(\mathbf{x}\right)\right\rangle + b\right) - y_{i} \leq \varepsilon + \xi_{i}^{*} \\ \xi_{i}, \xi_{i}^{*} \geq 0, \qquad i = 1, \dots, N \end{cases}$$

$$(2)$$

The constant C > 0 determines the trade-off between the smoothness of f and the cost associated to deviations larger than  $\varepsilon$  and  $\xi, \xi^*$  are the usual slack variables. The optimization problem (2) can be solved in dual form by introducing the positive Lagrange multipliers  $\alpha, \alpha^*, \eta, \eta^*$ . We can write the Lagrangian equation as follows:

$$L = \frac{1}{2} \|w\|^{2} + C \sum_{i=1}^{N} (\xi_{i} + \xi_{i}^{*}) - \sum_{i=1}^{N} (\eta_{i}\xi_{i} + \eta_{i}^{*}\xi_{i}^{*}) - \sum_{i=1}^{N} \alpha_{i} (\varepsilon + \xi_{i} + y_{i} - (\langle w, \Phi(\mathbf{x}_{i}) \rangle + b)) - \sum_{i=1}^{N} \alpha_{i}^{*} (\varepsilon + \xi_{i}^{*} - y_{i} + (\langle w, \Phi(\mathbf{x}) \rangle + b))$$
(3)

When minimizing L with respect to the primal variables  $(w, b, \xi_i, \xi_i^*)$ , the derivative has to vanish. It can be shown that the dual variables  $\eta_i, \eta_i^*$  are eliminated, so we obtain the dual optimization problem:

$$\max_{\alpha,\alpha^*} \quad W(\alpha^{(*)}) = -\frac{1}{2} \sum_{i,j=1}^N \left(\alpha_i - \alpha_i^*\right) \left(\alpha_j - \alpha_j^*\right) k(x_i, x_j) \\ -\varepsilon \sum_{i=1}^N \left(\alpha_i + \alpha_i^*\right) + \sum_{i=1}^N y_i \left(\alpha_i - \alpha_i^*\right)$$

subject to  $\begin{cases} \sum_{i=1}^{N} (\alpha_i - \alpha_i^*) = 0\\ \alpha_i, \alpha_i^* \in [0, C] \end{cases}$ 

where  $\alpha^{(*)} = (\alpha_1, \dots, \alpha_N, \alpha_1^*, \dots, \alpha_N^*)^T$  and  $k(x_i, x_j) = \langle \Phi(x_i), \Phi(x_j) \rangle$ . The estimate of the regression function at any given point x is then:

$$f(x) = \sum_{i=1}^{N} (\alpha_i - \alpha_i^*) k(x_i, x) + b$$
 (5)

where b is computed by exploiting the Karush-Kuhn-Tucker (KKT) conditions.

## 3. ALGORITHM DESCRIPTION

In this section, we present an overview of our proposed F-SVR learning algorithm for Support Vector Regression. The goal of the regression problem is to find the Lagrange multipliers  $\alpha^{(*)}$ , solution of the optimization problem (4). To solve this problem, we consider the particular analytical expression

of the linear constraint in (4). At least two components of  $\alpha$  and/or  $\alpha^*$  must be modified. We shall consider here the simplest case where only two components will be updated. Examining this linear constraint leads to four possible cases, where  $\lambda > 0$  is the only possibility to consider:

- case 1:  $\alpha_i \leftarrow \alpha_i + \lambda$ ;  $\alpha_j \leftarrow \alpha_j \lambda$
- case 2:  $\alpha_i^* \leftarrow \alpha_i^* + \lambda$ ;  $\alpha_i^* \leftarrow \alpha_i^* \lambda$
- case 3:  $\alpha_i \leftarrow \alpha_i + \lambda$ ;  $\alpha_i^* \leftarrow \alpha_i^* + \lambda$
- case 4:  $\alpha_i \leftarrow \alpha_i \lambda$ ;  $\alpha_i^* \leftarrow \alpha_i^* \lambda$

# 3.1. Searching for the *best* couple for each case (possible working set)

For every case, the modification of the selected components of  $\alpha$  and/or  $\alpha^*$  is defined by a vector  $u^{(*)} = (u_1, \ldots, u_N, u_1^*, \ldots, u_N^*)^T$  that has only two non-null components equal to  $\pm 1$ . In [5], it has been proposed to select the direction  $u^{(*)}$  in such a way that  $u^{(*)^T} g^{(*)}$  is maximum, where  $g^{(*)} = (g_1, \ldots, g_N, g_1^*, \ldots, g_N^*)^T$  is the gradient of the objective function  $W(\alpha^{(*)})$ . This is exactly what we use here:

• case 1: 
$$(i_1, j_1) = \underset{i,j}{\operatorname{arg\,max}} (g_i - g_j) \Leftrightarrow$$
  
 $i_1 = \underset{i}{\operatorname{arg\,max}} (g_i); j_1 = \underset{j}{\operatorname{arg\,min}} (g_j)$   
 $(u_{i1} = +1, u_{j1} = -1)$ 

- case 2:  $(i_2, j_2) = \operatorname*{arg\,max}_{i,j} \left(g_i^* g_j^*\right) \Leftrightarrow$   $i_2 = \operatorname*{arg\,max}_i(g_i^*); j_2 = \operatorname*{arg\,min}_j(g_j^*)$  $(u_{i2}^* = +1, u_{j2}^* = -1)$
- case 3:  $(i_3, j_3) = \underset{i,j}{\operatorname{arg\,max}} (g_i + g_j^*) \Leftrightarrow$   $i_3 = \underset{i}{\operatorname{arg\,max}} (g_i); j_3 = \underset{j}{\operatorname{arg\,max}} (g_j^*)$  $(u_{i3} = +1, u_{j3}^* = +1)$
- case 4:  $(i_4, j_4) = \operatorname*{arg\,max}_{i,j} (-g_i g_j^*) \Leftrightarrow$   $i_4 = \operatorname*{arg\,min}_i (g_i); j_4 = \operatorname*{arg\,min}_j (g_j^*)$  $(u_{i4} = -1, u_{j4}^* = -1)$

Before any modification of  $\alpha$  and/or  $\alpha^*$ , we must verify that the selected direction is feasible. A point  $\alpha^{(*)}$  is feasible if it satisfies all the constraints. A direction vector  $u^{(*)}$  is feasible, if there exists s > 0 such that  $\alpha^{(*)} + su^{(*)}$  is feasible. We must now take into account the inequality constraints. A direction will be feasible if the corresponding Lagrangian multipliers are not still located on the boundaries defined by the inequality constraints, so the search for the optimum directions defined previously becomes:

• case 1:  $i_1 = \underset{i:\alpha_i < C}{\operatorname{arg\,max}}(g_i); \ j_1 = \underset{j:\alpha_j > 0}{\operatorname{arg\,min}}(g_j)$ 

(4)

- case 2:  $i_2 = \underset{i:\alpha_i^* < C}{\operatorname{arg\,max}}(g_i^*); \ j_2 = \underset{j:\alpha_j^* > 0}{\operatorname{arg\,min}}(g_j^*)$
- case 3:  $i_3 = \underset{i:\alpha_i < C}{\operatorname{arg\,max}}(g_i); \ j_3 = \underset{j:\alpha_j^* < C}{\operatorname{arg\,max}}(g_j^*)$

• case 4: 
$$i_4 = \underset{i:\alpha_i>0}{\operatorname{arg\,min}}(g_i); \ j_4 = \underset{j:\alpha_j^*>0}{\operatorname{arg\,min}}(g_j^*)$$

### 3.2. Determination of the optimal step-size

Considering the quadratic nature of the objective function in (4),  $W(\alpha^{(*)})$ , it is actually easy to determine, for every case considered previously, the optimal step-size of the gradient algorithm. Given a direction  $u^{(*)}$ , a feasible point  $\alpha^{(*)}$  and in the unconstrained case, the optimal step-size  $\lambda_{opt}$  is defined by:  $\lambda_{opt} = \underset{\lambda}{\arg \max} \left( W(\alpha^{(*)} + \lambda u^{(*)}) \right)$ 

$$\lambda_{opt} = \frac{u^{(*)^T} \nabla W}{u^{(*)T} K u^{(*)}}$$
(6)

where K denotes the Gram matrix (that is not used, see below) and  $\nabla W$  is the gradient of  $W(\alpha^{(*)})$ . The optimal stepsize of equation (6) can be written as:

• case 1: 
$$\lambda_{opt_1} = \frac{g_{i_1} - g_{j_1}}{k_{i_1i_1} + k_{j_1j_1} - 2ki_1j_1}$$
  
• case 2:  $\lambda_{opt_2} = \frac{g_{i_2}^* - g_{j_2}^*}{k_{i_2i_2} + k_{j_2j_2} - 2ki_2j_2}$ 

• case 3: 
$$\lambda_{opt_3} = \frac{g_{i_3} + g_{j_3}^*}{k_{i_3i_3} + k_{j_3j_3} - 2ki_{3j_3}}$$

• case 4: 
$$\lambda_{opt_4} = \frac{-g_{i_4} - g_{j_4}^*}{k_{i_4i_4} + k_{j_4j_4} - 2ki_4j_4}$$

As we can see from the calculation, we are only working in dimension two. Now we have to verify that, after any modification, the updated point remains feasible. An analysis of the first case leads to:

$$\begin{cases} \alpha_{i_1} + \lambda_{opt_1} \leq C \\ \alpha_{j_1} - \lambda_{opt_1} \geq 0 \end{cases} \Rightarrow \lambda_{opt_1} \leq \min(C - \alpha_{i_1}, \alpha_{j_1})$$

If  $\lambda_{opt_1} > \min(C - \alpha_{i_1}, \alpha_{j_1})$  the retained value is equal to  $\min(C - \alpha_{i_1}, \alpha_{j_1})$ . A similar analysis for the other cases gives:

- case 2:  $\lambda_{opt_2} \leq \min(C \alpha_{i_2^*}, \alpha_{j_2^*})$ else  $\lambda_{opt_2} = \min(C - \alpha_{i_2^*}, \alpha_{j_2^*})$
- case 3:  $\lambda_{opt_3} \leq \min(C \alpha_{i_3}, C \alpha_{j_3^*})$ else  $\lambda_{opt_3} = \min(C - \alpha_{i_3}, C - \alpha_{j_3^*})$
- case 4:  $\lambda_{opt_4} \leq \min(\alpha_{i_4}, \alpha_{j_4^*})$ else  $\lambda_{opt_4} = \min(\alpha_{i_4}, \alpha_{j_4^*})$

## 3.3. Selection of the optimal case

We have now to select the best case in some sense. In our heuristic, that has been experimentally shown efficient, we select the case corresponding to the greatest increase of the objective function  $W(\alpha^{(*)})$ . Considering again the quadratic nature of  $W(\alpha^{(*)})$ , and after a few calculations, this increment is given by:

• case 1:  

$$\begin{aligned} \Delta W_1 &= -\lambda_{opt_1} \sum_{k=1}^{N} \left( k_{i_1k} - k_{j_1k} \right) \left( \alpha_k - \alpha_k^* \right) \\
&- \frac{1}{2} \lambda_{opt_1}^2 \left( k_{i_1i_1} + k_{j_1j_1} - 2k_{i_1j_1} \right) \\
&+ \lambda_{opt_1} \left( y_{i_1} - y_{j_1} \right) \end{aligned}$$
• case 2:  

$$\begin{aligned} \Delta W_2 &= -\lambda_{opt_2} \sum_{k=1}^{N} \left( k_{i_2k} - k_{j_2k} \right) \left( \alpha_k - \alpha_k^* \right) \\
&- \frac{1}{2} \lambda_{opt_1}^2 \left( k_{i_2i_2} + k_{j_2j_2} - 2k_{i_2j_2} \right) \\
&+ \lambda_{opt_1} \left( y_{i_2} - y_{j_2} \right) \end{aligned}$$
• case 3:  

$$\begin{aligned} \Delta W_3 &= -\lambda_{opt_3} \sum_{k=1}^{N} \left( k_{i_3k} - k_{j_3k} \right) \left( \alpha_k - \alpha_k^* \right) \\
&- \frac{1}{2} \lambda_{opt_3}^2 \left( k_{i_3i_3} + k_{j_3j_3} - 2k_{i_3j_3} \right) \\
&+ \lambda_{opt_3} \left( y_{i_3} - y_{j_3} \right) - 2\varepsilon \lambda_{opt_3} \end{aligned}$$
• case 4:  

$$\begin{aligned} \Delta W_4 &= -\lambda_{opt_4} \sum_{k=1}^{N} \left( k_{i_4k} - k_{j_4k} \right) \left( \alpha_k - \alpha_k^* \right) \\
&- \frac{1}{2} \lambda_{opt_4}^2 \left( k_{i_4i_4} + k_{j_4j_4} - 2k_{i_4j_4} \right) \\
&- \lambda_{opt_4} \left( y_{i_4} - y_{j_4} \right) + 2\varepsilon \lambda_{opt_4} \end{aligned}$$

The accepted case is the one that maximizes the increment of the objective function. We need to update the gradient of the objective function. The component l, (l = 1, ..., N) of the gradient is updated according to:

• case 1:  $\Delta g_l = -\lambda_{opt_1} \left( k_{i_1l} - k_{j_1l} \right)$  $\Delta g_l^* = \lambda_{opt_1} \left( k_{i_1l} - k_{j_1l} \right)$ 

• case 2: 
$$\begin{array}{l} \Delta g_l = \lambda_{opt_2} \left( k_{i_2l} - k_{j_2l} \right) \\ \Delta g_l^* = -\lambda_{opt_2} \left( k_{i_2l} - k_{j_2l} \right) \end{array}$$

• case 3: 
$$\begin{aligned} \Delta g_l &= -\lambda_{opt_3} \left( k_{i_3l} - k_{j_3l} \right) \\ \Delta g_l^* &= \lambda_{opt_3} \left( k_{i_3l} - k_{j_3l} \right) \end{aligned}$$

• case 4: 
$$\begin{aligned} \Delta g_l &= \lambda_{opt_4} \left( k_{i_4l} - k_{j_4l} \right) \\ \Delta g_l^* &= -\lambda_{opt_4} \left( k_{i_4l} - k_{j_4l} \right) \end{aligned}$$

#### 3.4. Initialization

The gradient of the objective function is given by:

$$g = \frac{\partial W(\alpha^{(*)})}{\partial \alpha} = -K(\alpha - \alpha^{*}) - \varepsilon \mathbf{1}_{N} + y$$

$$g^{*} = \frac{\partial W(\alpha^{(*)})}{\partial \alpha^{*}} = K(\alpha - \alpha^{*}) - \varepsilon \mathbf{1}_{N} - y$$
(7)

At the initialization we set the Lagrange multipliers  $\alpha^{(*)}$  to zero (feasible solution), then:

$$g = -\varepsilon 1_N + y$$
  
and  $W(0) = 0$  (8)  
$$g^* = -\varepsilon 1_N - y$$

## 3.5. Algorithm

F-SVR learning algorithm is resumed as follows:

- 1. Initialization.
- 2. Evaluation of the best direction, optimal step-size and increment of the objective function for each case.
- 3. Selection of the case that maximizes this increment.
- 4. Update of the solution, the gradient and the objective function.
- 5. Test for convergence, if not, return to step 2.

#### 4. EXPERIMENTAL RESULTS

We compare now our F-SVR algorithm with an implementation of SVMLight for regression problems, only for learning results comparison (no test data were used here). For SVM-Light, we consider different sizes of the working set. The first dataset is generated according to a noisy nonlinear model (see Fig 1). For this first experiment, both algorithms are written in MATLAB. Data are composed of 2000 observations of one predictor and the corresponding response. For this dataset,  $\varepsilon = 0.2, C = 10$  and the gaussian kernel parameter  $\sigma$  is equal to 0.3. Results are summarized in Table 1 where we evaluate the training Normalised Mean Square Error NMSE<sup>1</sup>.



Fig. 1. Simulated data (200 observations shown), SV are indicated by a circle

The second dataset (Census-house<sup>2</sup>) is concerned with predicting the median price of a house using some demographic information. Data are composed of 20000 observations (with 121 attributes). We compare our F-SVR algorithm to LIBSVM and SVMLIGHT<sup>3</sup> for regression problems. All algorithms are here written in C. For this experiment  $\varepsilon = 10^{-6}$ , C = 1000 and  $\sigma = 100$ . Obtained results are shown in Table 2.

	Time (sec)	Working set	NMSE
F-SVR algorithm	26.672	2	0.0048
SVMLight	199.125	2	0.024
	218.750	4	-
	253.016	8	-
	397.172	16	-
	722.328	32	-

Table 1. Results comparison for simulated data

	Time (sec)	Working set	NMSE
F-SVR algorithm	3742.99	2	0.0019
SVMLIGHT	5044.36	2	0.019
LIBSVM	5708.61	-	0.019

Table 2.	Results	comparison	for Census-	housing	data

#### 5. CONCLUSION

In this paper, we introduced a new learning algorithm for regression problems. This algorithm is easy to implement and does not require the storage of the Gram matrix. It efficiently uses the nature of the constraints of the optimization problem and properties of quadratic functions. We showed that this algorithm is fast and efficient for training large data sets, as proved our experiments. Development of an online version of our algorithm is under investigation.

#### 6. REFERENCES

- Smola A.J. and Scholkopf B., "A tutorial on support vector regression," *Statistics and Computing*, vol. 14, no. 3, pp. 199–222, 2004.
- [2] Platt J.C., "Fast training of support vector machines using sequential minimal optimization," Advances in Kernel Methods: Support Vector Learning, pp. 185–208, 1999.
- [3] S. Weston J. Bottou L. Bordes A., Ertekin, "Fast kernel classifiers with online and active learning.," *Journal of Machine Learning Research (JMLR)*, vol. 6, pp. 1579– 1619, 2005.
- [4] Bengio S. Collobert R., "Svmtorch: Support vector machines for large-scale regression problems," *Journal of Machine Learning Research*, no. 1, pp. 143–160, 2001.
- [5] Joachims T., "Making large-scale support vector machine learning practical," Advances in Kernel Methods: Support Vector Machines, pp. 169–184, 1999.
- [6] Vapnik V., *The nature of statistical learning theory*, Springer, New York, 1995.

 $<sup>^{1}</sup>$ NMSE =  $\frac{\sum_{i=1}^{N} (f(x_{i}) - y_{i})^{2}}{\sum_{i=1}^{N} (f(x_{i}) - y_{i})^{2}}$ 

<sup>&</sup>lt;sup>2</sup>http://www.cs.toronto.edu/ delve/data/census-house/

<sup>&</sup>lt;sup>3</sup>http://www.csie.ntu.edu.tw/ cjlin/libsvm/, http://svmlight.joachims.org/