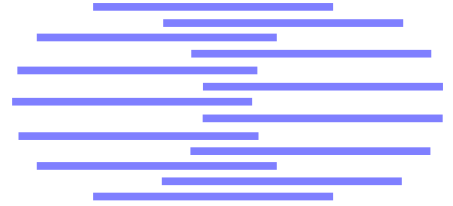


IDIAP

Martigny - Valais - Suisse



SUPPORT VECTOR MACHINES FOR LARGE-SCALE REGRESSION PROBLEMS

Ronan Collobert ¹ Samy Bengio ²

IDIAP-RR 00-17

AUGUST 16, 2000

Dalle Molle Institute
for Perceptual Artificial
Intelligence • P.O.Box 592 •
Martigny • Valais • Switzerland

phone +41 - 27 - 721 77 11
fax +41 - 27 - 721 77 12
e-mail secretariat@idiap.ch
internet <http://www.idiap.ch>

¹ IDIAP, CP 592, 1920 Martigny, Switzerland, collober@idiap.ch

² IDIAP, CP 592, 1920 Martigny, Switzerland, bengio@idiap.ch

SUPPORT VECTOR MACHINES FOR LARGE-SCALE REGRESSION PROBLEMS

Ronan Collobert

Samy Bengio

AUGUST 16, 2000

Abstract. Support Vector Machines (SVMs) for regression problems are trained by solving a quadratic optimization problem which needs on the order of l^2 memory and time resources to solve, where l is the number of training examples. In this paper, we propose a decomposition algorithm, *SVM Torch*¹, which is similar to *SVM-Light* proposed by Joachims [5] for classification problems, but adapted to regression problems. With this algorithm, one can now efficiently solve large-scale regression problems (more than 20000 examples). Comparisons with *Nodelib*, another SVM algorithm for large-scale regression problems from Flake and Lawrence [3] yielded significant time improvements.

¹*SVM Torch* is available at <http://www.idiap.ch/learning/SVMTorch.html>.

1 Introduction

Vapnik has proposed in [14] a method to solve regression problems using Support Vector Machines. It has yielded excellent performances on many regression and time series prediction problems (see for instance [9, 2]). This paper proposes an efficient implementation of SVMs for large-scale regression problems. Let us first recall how it works.

Given a training set of l examples (\mathbf{x}_i, y_i) with $\mathbf{x}_i \in E$ and $y_i \in \mathbb{R}$, where E is an Euclidean space with a scalar product denoted (\cdot) , we want to estimate the following linear regression:

$$f(\mathbf{x}) = (\mathbf{w} \cdot \mathbf{x}) + b$$

(with $b \in \mathbb{R}$) with a precision ϵ . For this, we minimize

$$\frac{1}{2}\|\mathbf{w}\|^2 + C \sum_{i=1}^l |y_i - f(\mathbf{x}_i)|_\epsilon$$

where $\frac{1}{2}\|\mathbf{w}\|^2$ is a regularization factor, C is a fixed constant, and $|\cdot|_\epsilon$ is the ϵ -insensitive loss function defined by Vapnik:

$$|z|_\epsilon = \max\{0, |z| - \epsilon\}.$$

Written as a constrained optimization problem, it amounts to minimizing

$$\tau(\mathbf{w}, \boldsymbol{\xi}, \boldsymbol{\xi}^*) = \frac{1}{2}\|\mathbf{w}\|^2 + C \sum_{i=1}^l (\xi_i + \xi_i^*)$$

subject to

$$((\mathbf{w} \cdot \mathbf{x}_i) + b) - y_i \leq \epsilon + \xi_i \tag{1}$$

$$y_i - ((\mathbf{w} \cdot \mathbf{x}_i) + b) \leq \epsilon + \xi_i^* \tag{2}$$

$$\xi_i, \xi_i^* \geq 0.$$

To generalize to non-linear regression, we replace the dot product with a kernel $k(\cdot)$. Then, introducing Lagrange multipliers $\boldsymbol{\alpha}$ and $\boldsymbol{\alpha}^*$, the optimization problem can be stated as:

Minimize the function

$$\mathcal{W}(\boldsymbol{\alpha}, \boldsymbol{\alpha}^*) = \frac{1}{2}(\boldsymbol{\alpha}^* - \boldsymbol{\alpha})^T K (\boldsymbol{\alpha}^* - \boldsymbol{\alpha}) - (\boldsymbol{\alpha}^* - \boldsymbol{\alpha})^T \mathbf{y} + \epsilon(\boldsymbol{\alpha}^* + \boldsymbol{\alpha})^T \mathbf{1} \tag{3}$$

subject to

$$(\boldsymbol{\alpha} - \boldsymbol{\alpha}^*)^T \mathbf{1} = 0 \tag{4}$$

and

$$0 \leq \alpha_i^*, \alpha_i \leq C, \quad i = 1 \dots l \tag{5}$$

where K is the matrix with coefficients $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$. The estimate of the regression function at a given point is then

$$f(\mathbf{x}) = \sum_{i=1}^l (\alpha_i^* - \alpha_i) k(\mathbf{x}_i, \mathbf{x}) + b$$

where b is computed using the fact that (1) becomes an equality with $\xi_i = 0$ if $0 < \alpha_i < C$ and (2) becomes an equality with $\xi_i^* = 0$ if $0 < \alpha_i^* < C$.

Solving the minimization problem (3) under the constraints (4) and (5) needs resources on the order of l^2 and is thus difficult for problems with large l .

In this paper, we propose a method to solve such problems efficiently using a decomposition algorithm similar to the one proposed by Joachims [5] in the context of classification problems. In the next section, we give the general algorithm and explain in more details each of its main steps, as well as a discussion on some important implementation issues, such as a way to efficiently handle the kernel matrix computation. In the experiment section, we first compare this new algorithm to *Nodelib* [3], another SVM algorithm for large-scale regression problems, and then show how the size of the internal memory allocated to the resolution of the problem is related to the time needed to solve it.

2 The Decomposition Algorithm

As in the classification algorithm proposed by Joachims [5], which was based on an idea from Osuna *et al* [10], our regression algorithm is subdivided into the following four steps:

1. Select q variables as the new working set, called \mathcal{S} .
2. Fix the other variables \mathcal{F} to their current values and solve the problem (3) with respect to \mathcal{S} .
3. Search for variables that are stuck to 0 or C and that will probably not change anymore. This is the *shrinking* phase.
4. Test if the optimization is finished or if we go back to the first step.

2.1 Selection of a New Working Set

We propose to select a new set of variables such that the overall criterion will be the most optimized. In order to select such working set, we use the same idea as Joachims [5]: simply search for the optimal gradient descent direction \mathbf{p} which is *feasible* and which has only q non-null components. The variables corresponding to these components will be our new working set \mathcal{S} .

We thus need to minimize:

$$\mathcal{V}(\mathbf{p}) = \left(\mathcal{W}'(\boldsymbol{\alpha}, \boldsymbol{\alpha}^*) \right)^T \mathbf{p} \quad (6)$$

with

$$\mathbf{p} = (d_1 \dots d_l, d_1^* \dots d_l^*)^T$$

subject to:

$$\mathbf{1}^T \mathbf{d} - \mathbf{1}^T \mathbf{d}^* = 0 \quad (7)$$

$$\begin{aligned} d_i &\geq 0 && \text{for } i \text{ such that } \alpha_i = 0 \\ d_i^* &\geq 0 && \text{for } i \text{ such that } \alpha_i^* = 0 \\ d_i &\leq 0 && \text{for } i \text{ such that } \alpha_i = C \\ d_i^* &\leq 0 && \text{for } i \text{ such that } \alpha_i^* = C \end{aligned} \quad (8)$$

and

$$-\mathbf{1} \leq \mathbf{p} \leq \mathbf{1} \quad (9)$$

$$\text{card}\{d_i : d_i \neq 0\} = q. \quad (10)$$

Since we are searching for an optimal descent direction, which is a direction where the scalar product with the gradient is the smallest, we want indeed to minimize (6). The conditions (7)-(8) are necessary to ensure the feasibility of the obtained direction. The condition (9) is there only to ensure that the problem has a solution. Finally, (10) is imposed because we are searching for a direction with only q non-null components.

Note that the derivative of \mathcal{W} can be easily computed:

$$\mathcal{W}'(\boldsymbol{\alpha}, \boldsymbol{\alpha}^*) = \begin{pmatrix} K(\boldsymbol{\alpha} - \boldsymbol{\alpha}^*) + \mathbf{y} + \mathbf{1} \epsilon \\ K(\boldsymbol{\alpha}^* - \boldsymbol{\alpha}) - \mathbf{y} + \mathbf{1} \epsilon \end{pmatrix}.$$

In order to solve this problem, it thus suffices to consider

$$\omega_i = \delta_i \mathcal{W}'_i$$

where $\delta_i = 1$ for $1 \leq i \leq l$ and $\delta_i = -1$ for $l+1 \leq i \leq 2l$. Let us force q to be even, and let us sort the ω_i in reverse order. Let us then denote φ as the bijection of $\{1 \dots 2l\}$ into itself such that the $\omega_{\varphi(i)}$ are sorted. Let us then select the $q/2$ first indices $\varphi(i)$ such that:

$$\text{if } \varphi(i) \leq l, \text{ we have } 0 < \alpha_{\varphi(i)} \leq C$$

$$\text{if } \varphi(i) > l, \text{ we have } 0 \leq \alpha_{\varphi(i)-l}^* < C$$

and let us select also the $q/2$ last indices $\varphi(i)$ such that:

$$\text{if } \varphi(i) \leq l, \text{ we have } 0 \leq \alpha_{\varphi(i)} < C$$

$$\text{if } \varphi(i) > l, \text{ we have } 0 < \alpha_{\varphi(i)-l}^* \leq C.$$

Since we are searching for exactly q variables, the $\varphi(i)$ must be distinct. We could have to reduce q if one variable is selected twice.

Let us now denote c_i , $i = 1 \dots q$ the q indices we just chose, we note that the direction which has for j^{th} component 0 if $j \notin \{c_1 \dots c_q\}$, $-\delta_j$ if $j \in \{c_1 \dots c_{\frac{q}{2}}\}$ and δ_j if $j \in \{c_{\frac{q}{2}+1} \dots c_q\}$, is a solution of the minimization problem (6)¹.

Our new working set \mathcal{S} is then composed of the q variables corresponding to the indices c_i (where an index $\leq l$ corresponds to α_{c_i} and an index $> l$ corresponds to $\alpha_{c_i-l}^*$).

¹To see that, let us go back to the minimization problem of

$$(z_1, \dots, z_l) \mapsto \sum_{i=1 \dots l} \omega_i z_i \tag{11}$$

subject to

$$\sum_i z_i = 0 \tag{12}$$

$$-1 \leq z_i \leq 1 \tag{13}$$

and

$$\text{card}\{z_i, z_i \neq 0\} = q \tag{14}$$

with $z_i = \delta_i p_i$. (The reasoning is the same if we take the constraints (8) into account).

In the case where $l = q = 2r$, it is easy to see that the minimum is obtained for $z_i = -1$ if $i = 1 \dots r$ and $z_i = 1$ if $i = r+1 \dots q$: if for instance z_{i_0} is augmented by $\gamma \geq 0$ for a $i_0 \leq r$, then one needs to compensate by $-\gamma$ another z_j to keep (12). Since we want to minimize (11), the best thing to do, knowing that the ω_i are sorted in reverse order and keeping in mind the constraint (13), is to modify z_q and thus to fix $z_q = 1 - \gamma$. Equation (11) is then augmented by $(\omega_{i_0} - \omega_q)\gamma$, which is a positive value because the ω_i are sorted and thus we get out of the minimum.

In the case where $l > q = 2r$, suppose we found a \mathbf{z} which is a solution of (11). Let us denote $k_1 \dots k_q$ the q indices of the components of \mathbf{z} which are non-nulls. Using the same argument as in the previous paragraph, it is clear that $z_{k_1} \dots z_{k_r} = -1$ and that $z_{k_{r+1}} \dots z_{k_q} = 1$. In other words, if \mathbf{z} is a solution to our problem, then we necessarily have $z_{k_i} = \pm 1$. Considering again the order of the ω_i , it becomes evident that we have to take $(k_1 = 1) \dots (k_r = r)$ and $(k_{r+1} = l - r) \dots (k_q = l)$.

2.2 Solving the Sub-Problem

We want to solve the problem (3) taking into account only variables \mathcal{S} . To simplify the notation, let us define

$$\boldsymbol{\beta} = \begin{pmatrix} \boldsymbol{\alpha} \\ -\boldsymbol{\alpha}^* \end{pmatrix}$$

and

$$\tilde{K} = \begin{pmatrix} K & K \\ K & K \end{pmatrix}$$

as well as

$$\mathbf{b} = \begin{pmatrix} -\mathbf{y} - \mathbf{1} \epsilon \\ -\mathbf{y} + \mathbf{1} \epsilon \end{pmatrix}.$$

The problem (3) is thus equivalent to minimize

$$\tilde{\mathcal{W}}(\boldsymbol{\beta}) = \frac{1}{2} \boldsymbol{\beta}^T \tilde{K} \boldsymbol{\beta} - \boldsymbol{\beta}^T \mathbf{b} \quad (15)$$

subject to:

$$\boldsymbol{\beta}^T \mathbf{1} = 0 \quad (16)$$

and

$$0 \leq \delta_i \beta_i \leq C, \quad i = 1 \dots 2l \quad (17)$$

where again $\delta_i = 1$ for $1 \leq i \leq l$ and $\delta_i = -1$ for $l+1 \leq i \leq 2l$.

Now let us suppose we can decompose each of the following variables into two parts (after having reordered the variables accordingly): the first one corresponds to variables \mathcal{S} and the second part corresponds to the fixed variables \mathcal{F} :

$$\boldsymbol{\beta} = \begin{pmatrix} \boldsymbol{\beta}_{\mathcal{S}} \\ \boldsymbol{\beta}_{\mathcal{F}} \end{pmatrix}$$

$$\mathbf{b} = \begin{pmatrix} \mathbf{b}_{\mathcal{S}} \\ \mathbf{b}_{\mathcal{F}} \end{pmatrix}$$

and

$$\tilde{K} = \begin{pmatrix} \tilde{K}_{\mathcal{S}\mathcal{S}} & \tilde{K}_{\mathcal{S}\mathcal{F}} \\ \tilde{K}_{\mathcal{F}\mathcal{S}} & \tilde{K}_{\mathcal{F}\mathcal{F}} \end{pmatrix}.$$

Replacing these variables in (15), (16) and (17), and taking into account the fact that $\tilde{K}_{\mathcal{S}\mathcal{F}}^T = \tilde{K}_{\mathcal{F}\mathcal{S}}$, the minimization problem is now

$$\tilde{\mathcal{W}}(\boldsymbol{\beta}_{\mathcal{S}}) = \frac{1}{2} \boldsymbol{\beta}_{\mathcal{S}}^T \tilde{K} \boldsymbol{\beta}_{\mathcal{S}} - \boldsymbol{\beta}_{\mathcal{S}}^T (\mathbf{b}_{\mathcal{S}} - \tilde{K}_{\mathcal{S}\mathcal{F}} \boldsymbol{\beta}_{\mathcal{F}}) \quad (18)$$

(removing the constants that depend only on \mathcal{F}), subject to

$$\boldsymbol{\beta}_{\mathcal{S}}^T \mathbf{1} = -\boldsymbol{\beta}_{\mathcal{F}}^T \mathbf{1} \quad (19)$$

and

$$0 \leq \tilde{\delta}_i \beta_{\mathcal{S}_i} \leq C, \quad i = 1 \dots q \quad (20)$$

where $\tilde{\delta}_i = 1$ if the i^{th} variable in the set \mathcal{S} corresponds to an α_i , $\tilde{\delta}_i = -1$ if it corresponds to an α_i^* .

Minimizing (18) under the constraints (19) and (20) can be realized using a constrained quadratic optimizer, such as a conjugate gradient method with projection or an interior point method [4].

Moreover, following Platt's idea in *SMO* [11], if one fixes the size of the working set \mathcal{S} to two, the problem can also be solved analytically.

This particular case is important because experimental results show that it often gives the fastest convergence times. We thus detailed it here. Let us simplify again the notation:

$$\begin{aligned}\boldsymbol{\beta}_{\mathcal{S}} &= \begin{pmatrix} z_1 \\ z_2 \end{pmatrix} \\ \mathbf{h} &= \mathbf{b}_{\mathcal{S}} - \tilde{K}_{\mathcal{S}\mathcal{F}}\boldsymbol{\beta}_{\mathcal{F}} \\ \zeta &= -\boldsymbol{\beta}_{\mathcal{F}}^{\top} \mathbf{1} \\ \tilde{K}_{\mathcal{S}} &= \begin{pmatrix} k_{11} & k_{12} \\ k_{21} & k_{22} \end{pmatrix}.\end{aligned}$$

Minimizing (18) under the constraints (19) and (20) is thus equivalent to minimizing

$$(z_1, z_2) \mapsto \frac{1}{2} (k_{11} z_1^2 + k_{22} z_2^2 + 2k_{12} z_1 z_2) - h_1 z_1 - h_2 z_2 \quad (21)$$

subject to

$$z_1 + z_2 = \zeta \quad (22)$$

and

$$0 \leq \tilde{\delta}_1 z_1, \tilde{\delta}_2 z_2 \leq C. \quad (23)$$

We are searching for a minimum in (21) with respect to z_1 along the line (22). By inserting (22) into (21), and after some derivations, it is now equivalent to minimizing

$$\Lambda : z_1 \mapsto \frac{1}{2} (k_{11} - 2k_{12} + k_{22}) z_1^2 + [(k_{12} - k_{22}) \zeta - h_1 + h_2] z_1.$$

In the case² where $\eta = k_{11} - 2k_{12} + k_{22} > 0$, this function has a unique minimum for

$$z_1^o = \frac{(k_{22} - k_{12}) \zeta + h_1 - h_2}{\eta}.$$

Let us now consider the constraints (22) and (23). They force z_1 to stay between L and H where

$$\begin{aligned} \left. \begin{aligned} L &= \max(0, \zeta - C) \\ H &= \min(C, \zeta) \end{aligned} \right\} && \text{if } \tilde{\delta}_1 = 1 \text{ and } \tilde{\delta}_2 = 1 \\ \left. \begin{aligned} L &= \max(0, \zeta) \\ H &= \min(C, \zeta + C) \end{aligned} \right\} && \text{if } \tilde{\delta}_1 = 1 \text{ and } \tilde{\delta}_2 = -1 \\ \left. \begin{aligned} L &= \max(-C, \zeta - C) \\ H &= \min(0, \zeta) \end{aligned} \right\} && \text{if } \tilde{\delta}_1 = -1 \text{ and } \tilde{\delta}_2 = 1 \\ \left. \begin{aligned} L &= \max(-C, \zeta) \\ H &= \min(0, \zeta + C) \end{aligned} \right\} && \text{if } \tilde{\delta}_1 = -1 \text{ and } \tilde{\delta}_2 = -1. \end{aligned}$$

Thus, taking

$$z_1^{o,c} = \begin{cases} H & \text{if } z_1^o > H \\ z_1^o & \text{if } L \leq z_1^o \leq H \\ L & \text{if } z_1^o < L \end{cases}$$

²Note that this is the most common case. For instance for a gaussian kernel with distinct examples \mathbf{x}_i , it is easy to see that is is always the case.

and

$$z_2^o = \zeta - z_1^{o,c}$$

the minimum of (21) under the constraints (22) and (23) is obtained at $(z_1^{o,c}, z_2^o)$ if $\eta > 0$. In the pathological case where $\eta \leq 0$, it is clear that the solution

$$z_1^o = \begin{cases} L & \text{if } \Lambda(L) < \Lambda(H) \\ H & \text{if } \Lambda(L) \geq \Lambda(H) \end{cases}$$

and

$$z_2^o = \zeta - z_1^o$$

is the minimum.

2.3 Shrinking

The idea of shrinking is to remove some variables that are stuck to the bounds 0 or C , and that will *a priori* not change anymore. To do this, we use the fact that (α, α^*) minimizes the problem (3) under the constraints (4) and (5) *if and only if* there exists numbers $\lambda^{up} \in \mathbb{R}^{2l}$, $\lambda^{low} \in \mathbb{R}^{2l}$, $\lambda^{eq} \in \mathbb{R}$ that verify the following *KKT conditions*:

$$\mathcal{W}'(\alpha, \alpha^*) + \lambda^{eq} \begin{pmatrix} 1 \\ -1 \end{pmatrix} - \lambda^{low} + \lambda^{up} = \mathbf{0} \quad (24)$$

$$\lambda_i^{low} \alpha_i = 0, \quad i = 1 \dots l \quad \text{and} \quad \lambda_i^{low} \alpha_{i-l}^* = 0, \quad i = (l+1) \dots 2l \quad (25)$$

$$\lambda_i^{up} (\alpha_i - C) = 0, \quad i = 1 \dots l \quad \text{and} \quad \lambda_i^{up} (\alpha_{i-l}^* - C) = 0, \quad i = (l+1) \dots 2l \quad (26)$$

$$\lambda^{low} \geq \mathbf{0} \quad (27)$$

$$\lambda^{up} \geq \mathbf{0} \quad (28)$$

$$(\alpha - \alpha^*)^T \mathbf{1} = 0 \quad (29)$$

$$\mathbf{0} \leq \alpha^*, \quad \alpha \leq C. \quad (30)$$

Note that if $\lambda_i^{low} > 0$, then the corresponding variable is equal to 0. Also, if $\lambda_i^{up} > 0$, then the corresponding variable is equal to C . The idea is thus to search at each iteration for variables λ^{up} , λ^{low} and λ^{eq} that verify *as well as possible*³ the equations (24)-(28), and to remove a variable which is stuck at 0 if its coefficient λ_i^{low} is strictly positive (or just below a constant ϵ_{shrink}) during a given number of iterations. Using the same idea, we also eliminate a variable which is stuck at C if its coefficient λ_i^{up} stays strictly positive for a sufficient number of iterations.

To estimate λ^{low} or λ^{up} , we start by estimating λ^{eq} (note that if $0 < \alpha_i < C$ then $\lambda_i^{low} = \lambda_i^{up} = 0$ and if $0 < \alpha_i^* < C$ then $\lambda_{i+l}^{low} = \lambda_{i+l}^{up} = 0$). Noting

$$A = \{i, 0 < \alpha_i < C\}, \quad B = \{i, 0 < \alpha_i^* < C\}$$

we have (with $\hat{\lambda}$ standing for an estimation of λ) :

$$\hat{\lambda}^{eq} = \frac{1}{|A \cup B|} \left(\sum_{i \in B} \mathcal{W}'_{i+l}(\alpha, \alpha^*) - \sum_{i \in A} \mathcal{W}'_i(\alpha, \alpha^*) \right). \quad (31)$$

³If we were *really* able to find such variables, this would mean that (α, α^*) is a solution to our problem.

Then if we have

$$\begin{array}{ll}
\alpha_i = 0 & \text{we compute } \hat{\lambda}_i^{low} = \hat{\lambda}^{eq} + \mathcal{W}'_i \\
\alpha_i^* = 0 & \text{we compute } \hat{\lambda}_{i+l}^{low} = -\hat{\lambda}^{eq} + \mathcal{W}'_{i+l} \\
\alpha_i = C & \text{we compute } \hat{\lambda}_i^{up} = -\hat{\lambda}^{eq} - \mathcal{W}'_i \\
\alpha_i^* = C & \text{we compute } \hat{\lambda}_{i+l}^{up} = \hat{\lambda}^{eq} - \mathcal{W}'_{i+l}
\end{array} \tag{32}$$

and if a variable stays a sufficient number of iterations at 0 (or C) with its corresponding coefficient $\hat{\lambda}_j^{low} > \epsilon_{shrink}$ (or $\hat{\lambda}_j^{up} > \epsilon_{shrink}$), then we remove it from the problem.

2.4 Termination Criterion

Given what has been said in the section on *shrinking*, if we can always have (29) and (30) during the resolution of a sub-problem, a *reasonable* termination criterion is to verify that the λ estimated by (31) and (32) verifies the conditions (24)-(28) with a given precision ϵ_{end} .

Thus we simply verify that

$$\begin{array}{ll}
\text{for } i \text{ such that } 0 < \delta_i \beta_i < C : & \lambda^{eq} - \epsilon_{end} \leq -\delta_i \mathcal{W}'_i(\boldsymbol{\alpha}, \boldsymbol{\alpha}^*) \leq \lambda^{eq} + \epsilon_{end} \\
\text{for } i \text{ such that } \beta_i = 0 : & \mathcal{W}'_i(\boldsymbol{\alpha}, \boldsymbol{\alpha}^*) + \delta_i \lambda^{eq} \geq -\epsilon_{end} \\
\text{for } i \text{ such that } \beta_i = C : & \mathcal{W}'_i(\boldsymbol{\alpha}, \boldsymbol{\alpha}^*) + \delta_i \lambda^{eq} \leq \epsilon_{end}
\end{array}$$

with $\delta_i = 1$ for $1 \leq i \leq l$, $\delta_i = -1$ for $(l+1) \leq i \leq 2l$ and

$$\boldsymbol{\beta} = \begin{pmatrix} \boldsymbol{\alpha} \\ -\boldsymbol{\alpha}^* \end{pmatrix}.$$

2.5 Implementation Details

Note that in all the steps needed during one iteration, only two of them could be time consuming: the one that computes \mathcal{W}'_i and the one that computes $\mathbf{b}_S - \tilde{K}_{S\mathcal{F}}\boldsymbol{\beta}_{\mathcal{F}}$ in (18).

We therefore propose to keep in memory a table of \mathcal{W}'_i . Moreover, to update this variable, we can see that

$$\begin{array}{ll}
\text{for } i \leq n & \mathcal{W}'_i^{(t+1)} = \mathcal{W}'_i^{(t)} + \sum_{j \in S_1} K_{ij} \left(\alpha_j^{(t+1)} - \alpha_j^{(t)} \right) - \sum_{j \in S_2} K_{ij} \left(\alpha_j^{*(t+1)} - \alpha_j^{*(t)} \right) \\
\text{for } i > n & \mathcal{W}'_i^{(t+1)} = \mathcal{W}'_i^{(t)} - \sum_{j \in S_1} K_{ij} \left(\alpha_j^{(t+1)} - \alpha_j^{(t)} \right) + \sum_{j \in S_2} K_{ij} \left(\alpha_j^{*(t+1)} - \alpha_j^{*(t)} \right)
\end{array}$$

where

$$S_1 = \{i, \alpha_i \in \mathcal{S}\} \text{ and } S_2 = \{i, \alpha_i^* \in \mathcal{S}\}.$$

For the computation of $\mathbf{b}_S - \tilde{K}_{S\mathcal{F}}\boldsymbol{\beta}_{\mathcal{F}}$, we can use the following trick:

$$\begin{aligned}
\left(\mathbf{b}_S - \tilde{K}_{S\mathcal{F}}\boldsymbol{\beta}_{\mathcal{F}} \right)_i &= \left(\mathbf{b}_S \right)_i - \left(\tilde{K}_{S\mathcal{F}}\boldsymbol{\beta}_{\mathcal{F}} + \tilde{K}_{SS}\boldsymbol{\beta}_S \right)_i + \left(\tilde{K}_{SS}\boldsymbol{\beta}_S \right)_i \\
&= \begin{cases} -\mathcal{W}_i + \left(\tilde{K}_{SS}\boldsymbol{\beta}_S \right)_i & \text{if } i \leq n \\ \mathcal{W}_i + \left(\tilde{K}_{SS}\boldsymbol{\beta}_S \right)_i & \text{if } i > n. \end{cases}
\end{aligned}$$

With these two ideas, one can reduce considerably the computational time: instead of computing all the lines of the matrix K , one can only compute the lines corresponding to the variables in \mathcal{S} .

Since we only need these lines for the computations, and since it quickly becomes intractable for large problems to keep all the matrix K in memory (the size of the matrix being quadratic with respect to the number of examples), it is interesting to implement a *cache* that keeps in memory the lines of K that corresponds to the most used variables instead of recomputing them at each iteration.

3 Experimental Results

3.1 Speed Comparisons

We compared our SVM implementation for regression problems (named *SVMTorch*) to the one from Flake and Lawrence [3] using their publicly available software *Nodelib*, which is an enhanced version of *SMO* [11].

Both these algorithms use an internal cache in order to be able to solve large-scale problems. All the experiments presented in this section have been done on a SPARC Ultra-10 440Mhz, with the *gcc* compiler. The parameters of the algorithms were not chosen to obtain the best generalization performances, since the goal was to *compare the speed* of the algorithms. Both programs used the same parameters with regard to cache, precision, etc... For *Nodelib*, the other parameters were set using the default values proposed by the authors. All the programs were compiled using *double* precision. Finally, each program was trained/tested on the same data for a given *benchmark*.

We compared the programs on two different tasks, using up to three different training set sizes. The first task was to predict the average number of sunspots of one year, given the average number of sunspots of the previous 12 years. Since we had access to daily data, we were able to artificially create one input/output pair for each day: at each day, we compute the yearly average of the year starting the next day, which has to be predicted using the 12 previous yearly averages. Using this technique, we had access to 43000 examples with input dimension equal to 12. The last 3000 examples served as the test set, while we created 3 different training sets of varying sizes: 5000, 20000 and 40000. The second task was created artificially using the daily sunspot series in order to test the algorithm with a big input dimension. We selected it to be equal to 240. The test set had 2000 examples while we created 2 different training sets of varying sizes: 5000 and 20000.

Table 1 gives the computation time results of all the experiments in CPU-seconds. In these experiments, the following SVM parameters were used: $C = 1000$ and $\text{precision} = 0.01$. As it can be seen, *SVMTorch* easily outperformed *Nodelib* in all experiments: the generalization performance of both algorithms as well as the number of support vectors found were similar but the time spent to find the solution was largely less for *SVMTorch* than for *Nodelib*. However, it seems that as the number of training examples grows, the speed difference between both algorithms narrows. The speed difference is also smaller when the input dimension is higher. A further analysis should probably be done in order to understand these difference variations.

Yearly Sunspots ($\sigma = 900, \epsilon = 20$)					
	Train size	# of SV	Train MSE	Test MSE	# of CPU-seconds
Nodelib	5000	431	146.15	246.18	245
SVMTorch		432	146.14	246.15	3
Nodelib	20000	1151	126.59	287.60	1253
SVMTorch		1155	126.58	287.64	30
Nodelib	40000	1911	119.46	363.88	2499
SVMTorch		1871	119.76	369.22	94
Artificial Data ($\sigma = 10, \epsilon = 0.5$)					
	Train size	# of SV	Train MSE	Test MSE	# of CPU-seconds
Nodelib	5000	397	0.09	4.82	825
SVMTorch		394	0.10	4.66	90
Nodelib	20000	1013	0.08	0.60	5598
SVMTorch		1002	0.08	0.60	1037

Table 1: Comparative speed results between *SVMTorch* and *Nodelib* on two different databases and up to 3 different training set sizes. All the experiments used an internal cache of 100Mb.

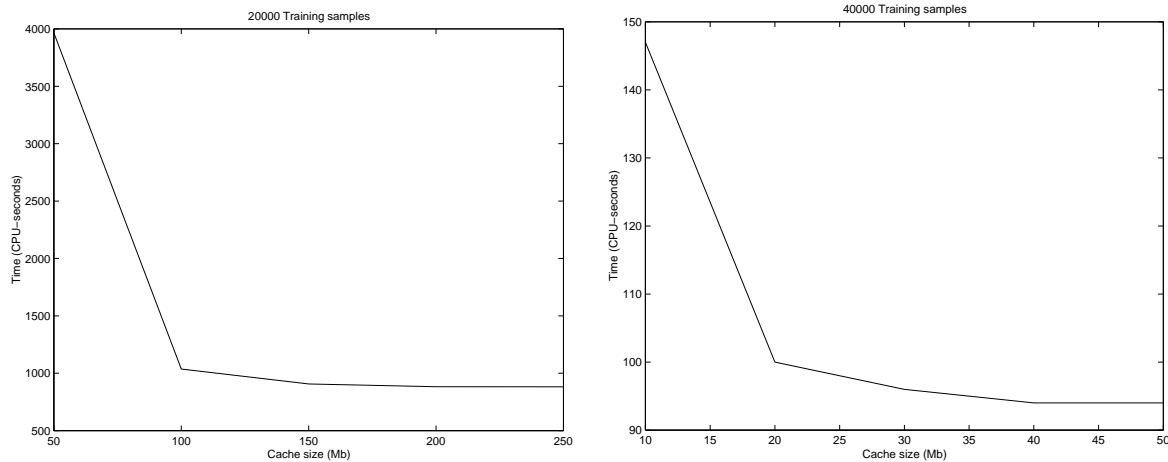


Figure 1: Evolution of the convergence speed (in seconds) with respect to the size of the internal cache (in Mb). The left part gives the computational time for the artificial dataset trained with 20000 examples (in dimension 240), while the right part gives the computational time for the sunspots dataset trained with 40000 examples (in dimension 12).

3.2 Analysis of the Cache Size

In order to show the impact of the cache size on the time needed to train an SVM, we have done two series of experiments, one for each dataset, varying the size of the cache (from 10 to 50 Mb for the sunspots dataset and from 50 to 250 for the artificial dataset). Figure 1 shows the results: the more cache you can afford, the faster the algorithm will find a solution, but the optimal size of the cache is of course related to the number of training samples: after a given cache size, no more speed improvement can be achieved.

4 Conclusion

We have presented a new decomposition algorithm intended to efficiently solve large-scale regression problems using SVMs. This algorithm followed the same principles as the one from Joachims [5] for his classification algorithm. We also proposed to set the size of the sub-problems to 2 in order to solve it analytically as it is done in SMO [11]. An internal cache keeping part of the kernel matrix in memory enables the program to solve large problems without the need to keep quadratic resources in memory and without the need to recompute every kernel evaluation, which leads to an overall fast algorithm. An experimental comparison with another algorithm has shown significant time improvement for large-scale problems.

Finally, note that in the available source code of *SVM-Torch*, we also implemented the same algorithm for classification problems. The result is thus similar mathematically to the one proposed by Joachims in its *SVM-Light* software [5], but the speed obtained was twice as fast as *SVM-Light* (there are probably many implementation differences, such as the cache, which was not fully described in Joachims's paper).

References

- [1] C. Chang, C. Hsu, and C. Lin. The analysis of decomposition methods for support vector machines. In *IJCAI'99, Workshop on Support Vector Machines, 1999*.

- [2] H. Drucker, C.J.C. Burges, L. Kaufman, A. Smola, and V. Vapnik. Support vector regression machines. In M. Mozer, M. Jordan, and T. Petsche, editors, *Advances in Neural Information Processing Systems 9*, pages 155–161. MIT Press, 1997.
- [3] G.W. Flake and S. Lawrence. Efficient SVM regression training with SMO. Submitted to Machine Learning. Available at <http://external.nj.nec.com/homepages/flake/smorch.ps>.
- [4] R. Fletcher. *Practical Methods of Optimization*. John Wiley and sons, 1987.
- [5] Thorsten Joachims. Making large-scale support vector machine learning practical. In B. Schölkopf, C. Burges, and A. Smola, editors, *Advances in Kernel Methods*. MIT Press, 1999.
- [6] S.S. Keerthi and E.G. Gilbert. Convergence of a generalized SMO algorithm for SVM classifier design. Technical Report CD-00-01, Control Division, Dept. of Mechanical and Production Engineering, National University of Singapore, 2000. Available at http://guppy.mpe.nus.edu.sg/~mpessk/svm/conv_ml.ps.gz.
- [7] S.S. Keerthi, S.K. Shevade, C. Bhattacharyya, and K.R.K. Murthy. Improvements to platt's SMO algorithm for SVM classifier design. Technical Report CD-99-14, Control Division, Dept. of Mechanical and Production Engineering, National University of Singapore, 1999. To appear in Neural Computation. Available at http://guppy.mpe.nus.edu.sg/~mpessk/smo_mod.ps.gz.
- [8] P. Laskov. An improved decomposition algorithm for regression support vector machines. In S.A. Solla, T.K. Leen, and K.-R. Müller, editors, *Advances in Neural Information Processing Systems 12*. MIT Press, 2000. Available at <http://www.cis.udel.edu/~laskov/publications/NIPS-99.ps.gz>.
- [9] K.-R. Müller, A. Smola, G. Rätsch, B. Schölkopf, J. Kohlmorgen, and V. Vapnik. Predicting time series with support vector machines. In W. Gerstner, A. Germond, M. Hasler, and J.-D. Nicoud, editors, *Artificial Neural Networks - ICANN'97*, pages 999–1004. Springer, 1997.
- [10] Edgar Osuna, Robert Freund, and Federico Girosi. An improved training algorithm for support vector machines. In J. Principe, L. Giles, N. Morgan, and E. Wilson, editors, *Neural Networks for Signal Processing VII - Proceedings of the 1997 IEEE Workshop*, pages 276–285. IEEE, New York, 1997.
- [11] John C. Platt. Fast training of support vector machines using sequential minimal optimization. In B. Schölkopf, C. Burges, and A. Smola, editors, *Advances in Kernel Methods*. MIT Press, 1999.
- [12] S.K. Shevade, S.S. Keerthi, C. Bhattacharyya, and K.R.K. Murthy. Improvements to smo algorithm for svm regression. Technical Report CD-99-16, Control Division, Dept. of Mechanical and Production Engineering, National University of Singapore, 1999. To appear in IEEE Transaction on Neural Networks. Available at http://guppy.mpe.nus.edu.sg/~mpessk/smoreg_mod.shtml.
- [13] A. J. Smola and B. Schölkopf. A tutorial on support vector regression. Technical Report NeuroCOLT Technical Report NC-TR-98-030, Royal Holloway College, University of London, UK, 1998.
- [14] Vladimir N. Vapnik. *The nature of statistical learning theory*. Springer, second edition, 1995.

A Things that should have been in the paper

A.1 On the experiments

We forgot to specify that the kernel used in all the experiments was a Gaussian kernel. Thanks to Flake⁴ for his remark on this.

In the conclusion of our paper, we said that we also implemented a classification version of the algorithm which was similar to the one proposed by Joachims [5] and that our implementation was twice as fast as *SVM-Light*. This assertion holds only for non-sparse data because *SVM-Light* has been specially designed for sparse data, while it was not the case in the first version of *SVM-Torch*. The current version, which now includes sparse data format, is 1.33 times faster than *SVM-Light* for sparse data (and still 2 times faster for non-sparse data).

A.2 On the decomposition method

Since the publication of our technical report, we have been aware of many other decomposition algorithms for regression problems that we have not even cited. We try here to resume their work and the relation with our paper.

Shevade *et al* [12] proposed two modifications of the *SMO* algorithm for regression, based on a previous paper from the same team [7] for the classification problem. Laskov [8] proposed also a decomposition method for regression problems which is very similar to the second one from Shevade *et al*. In fact, it is easy to see that Laskov's method with a subproblem of size 2 uses the same selection algorithm as well as the same termination criterion.

Their method for selecting the working set is *very* similar to the one we proposed, but while we propose to select variables α_i independantly of their counterpart α_i^* , they propose to select simultaneously *pairs* of variables $\{\alpha_i, \alpha_i^*\}$. Even if this seems to be a small difference, let us note that since $\alpha_i \alpha_i^* = 0 \forall i$, one of the two variables α_i or α_i^* is always equal to 0, and choosing the α_i and the α_i^* independantly can thus help to quickly eliminate half of the variables, thanks to the *shrinking* phase⁵, which of course have a direct impact on the speed of our program.

Similarly, Smola and Schölkopf [13] also proposed earlier to use a decomposition algorithm for regression based on *SMO* using an analytical solution for the subproblem, but again they propose to select 2 pairs of variables (2 α and their corresponding α^*) instead of 2 variables, which leads to a different mathematical formulation. As for Laskov and Shevada *et al*, they do not use *shrinking* which is, in our opinion, the main speed gain of our algorithm.

Finally, Flake and Lawrence [3] proposed again a modification of *SMO* for regression which uses the heuristics proposed by Platt [11] and those from Smola and Schölkopf [13] but works on a new variable $\lambda_i = \alpha_i - \alpha_i^*$, which leads to a different analytical solution. In fact, all the analytical solutions proposed by these authors are different but need to handle multiple cases for the solution, except our method.

A.3 On the convergence of the algorithm

In our paper, we did not talked about the convergence of our algorithm. A paper from Chang *et al* [1] talks about the convergence of some SVM algorithms based on a decomposition method. However, using their arguments, we cannot conclude that our algorithm converges to the optimum, even when no shrinking is done. The hypothesis they use in their proof regarding their method to search for a feasible solution is slightly different from our method and thus we cannot use their proof in our case. Keerthi *et al* [6] also proposed a convergence proof for their method [7], but it applies only to their classification case. However, we will see in the next section that it also applies to our classification method *as well as* our regression method.

⁴<http://www.neci.nj.nec.com/homepages/flake/>

⁵this is verified in practice

B Remarks on the relation between many SVM algorithms

As we said in our paper, the algorithm we used in classification is the same, mathematically speaking, as the one proposed by Joachims [5]. Let us now consider⁶ the paper from Keerthi *et al* [7] which proposes two algorithms based on Platt's algorithm, *SMD*. In particular, let us focus on the second method they propose and let us compare this method to the algorithm proposed by Joachims in the case of a working set of size 2.

We strongly suggest to the reader to refer to the papers from Joachims and Keerthi *et al* in order to understand the following notations which will not be re-explained here.

At each iteration, Keerthi *et al* propose to start by selecting two variables in their working set. Following their notation, let us denote

$$\begin{aligned} I_0 &= \{i : 0 < \alpha_i < C\} \\ I_1 &= \{i : y_i = 1, \alpha_i = 0\} \\ I_2 &= \{i : y_i = -1, \alpha_i = C\} \\ I_3 &= \{i : y_i = 1, \alpha_i = C\} \\ I_4 &= \{i : y_i = -1, \alpha_i = 0\} \end{aligned}$$

and let us denote also i_{low} and i_{up} the index of the two selected variables. They verify

$$F_{i_{low}} = b_{low} = \max\{F_i : i \in I_0 \cup I_3 \cup I_4\}$$

and

$$F_{i_{up}} = b_{up} = \min\{F_i : i \in I_0 \cup I_1 \cup I_2\}$$

where

$$F_i = \sum_j \alpha_j y_j k(\mathbf{x}_i, \mathbf{x}_j) - y_i.$$

One can easily remark that $F_i = \omega_i$, where ω_i is the sorting variable used by Joachims in his paper. Joachims defines the following constraints

$$\begin{aligned} d_i &\geq 0, \quad \forall i : \alpha_i = 0 \\ d_i &\leq 0, \quad \forall i : \alpha_i = C \end{aligned} \tag{33}$$

and select the following working set variables:

- the one that corresponds to the highest ω_i such that $0 < \alpha_i < C$ or such that $d_i = -y_i$ verify (33), and
- the one that corresponds to the smallest ω_i such that $0 < \alpha_i < C$ or such that $d_i = y_i$ verify (33).

This is indeed equivalent to the choice made by Keerthi *et al*.

Both algorithms then solve the sub-problem and test the optimality of the general problem. The algorithm from Keerthi *et al* stops when

$$b_{low} - b_{up} \leq \tau$$

where τ is a tolerance factor defined by the user. The algorithm from Joachims stops when the following conditions are verified:

$$\begin{aligned} \forall i \text{ such that } 0 < \alpha_i < C, \quad \lambda^{eq} - \tau \leq y_i - \sum_j \alpha_j y_j k(\mathbf{x}_i, \mathbf{x}_j) \leq \lambda^{eq} + \tau \\ \forall i \text{ such that } \alpha_i = 0, \quad y_i (\sum_j \alpha_j y_j k(\mathbf{x}_i, \mathbf{x}_j) + \lambda^{eq}) \geq 1 - \tau \\ \forall i \text{ such that } \alpha_i = C, \quad y_i (\sum_j \alpha_j y_j k(\mathbf{x}_i, \mathbf{x}_j) + \lambda^{eq}) \leq 1 + \tau \end{aligned} \tag{34}$$

⁶Thanks to Patrick Haffner (<http://www.research.att.com/~haffner>) and Ryan Rifkin (<http://five-percent-nation.mit.edu/PersonalPages/rif/>) who have stimulated our interest on this.

where λ^{eq} is defined as follows

$$\lambda^{eq} = \frac{1}{|A|} \sum_{i \in A} \left[y_i - \sum_j \alpha_j y_j k(\mathbf{x}_i, \mathbf{x}_j) \right]$$

with

$$A = \{i : 0 < \alpha_i < C\}.$$

It is easy to see that equations (34) are equivalent to

$$\begin{aligned} \forall i \in I_0, \quad \lambda^{eq} - \tau &\leq -F_i \leq \lambda^{eq} + \tau \\ \forall i \in I_1 \cup I_2, \quad F_i &\geq -\lambda^{eq} - \tau \\ \forall i \in I_3 \cup I_4, \quad F_i &\leq -\lambda^{eq} + \tau. \end{aligned}$$

Moreover, since $-b_{low} \leq \lambda^{eq} \leq -b_{up}$, if the optimality conditions from Keerthi *et al* are verified, then

$$\forall i \in I_0 \cup I_3 \cup I_4, \quad F_i \leq b_{up} + \tau \leq -\lambda^{eq} + \tau$$

and

$$\forall i \in I_0 \cup I_1 \cup I_2, \quad -F_i \leq -b_{low} + \tau \leq \lambda^{eq} + \tau$$

which implies the optimality conditions from Joachims.

Since the optimality test of *SVM-Light* is weaker than the one from Keerthi *et al* [6], it is easy to see that one can apply their theorem to show that *SVM-Light* converges for subproblems of size 2 (as well as our classification algorithm).

In the same way, one can show that the optimality test we used in our regression algorithm is weaker than the general algorithm proposed by Keerthi *et al* [6] and it is easy to see that given the fact that our algorithm uses a selection method that choose independantly the α and the α^* , the proof from Keerthi *et al* also applies to our regression algorithm when the subproblem size is set to 2.

C Conclusion

In conclusion, we note that all these decomposition algorithms are extremely related. For instance, subset selection algorithms from Keerthi *et al* and Joachims are strictly identical if the *shrinking* is not used and for a working subset of size 2 in *SVM-Light*.

Moreover, it is very easy to see that the regression method from Laskov (again with a subset of size 2) is equivalent to the one from Shevade *et al*, which is the same as the classification one from Keerthi *et al*. Note also that the method from Flake and Lawrence could be modified using the second modification from Keerthi *et al* and would thus be enhanced.

Finally, we think that shrinking makes the main difference with regard to speed, and the selection method we have chosen simplifies the resolution of the analytic quadratic problem and enables to obtain a convergence proof for the regression problem.