



# Feasible Direction Decomposition Algorithms for Training Support Vector Machines

PAVEL LASKOV

laskov@cis.udel.edu

Department of Computer and Information Sciences, 102 Smith Hall, University of Delaware, Newark,  
DE 19718, USA

**Editor:** Nello Cristianini

**Abstract.** The article presents a general view of a class of decomposition algorithms for training Support Vector Machines (SVM) which are motivated by the method of feasible directions. The first such algorithm for the pattern recognition SVM has been proposed in Joachims, T. (1999, Schölkopf et al. (Eds.) *Advances in kernel methods-Support vector learning* (pp. 185–208). MIT Press). Its extension to the regression SVM—the maximal inconsistency algorithm—has been recently presented by the author (Laskov, 2000, Solla, Leen, & Müller (Eds.) *Advances in neural information processing systems 12* (pp. 484–490). MIT Press). A detailed account of both algorithms is carried out, complemented by theoretical investigation of the relationship between the two algorithms. It is proved that the two algorithms are equivalent for the pattern recognition SVM, and the feasible direction interpretation of the maximal inconsistency algorithm is given for the regression SVM. The experimental results demonstrate an order of magnitude decrease of training time in comparison with training without decomposition, and, most importantly, provide experimental evidence of the linear convergence rate of the feasible direction decomposition algorithms.

**Keywords:** support vector machines, training, decomposition algorithms, methods of feasible directions, working set selection

## 1. Introduction

Computational complexity of SVM training algorithms is attracting increasing interest, as applications of SVM extend to problems of larger and larger size. It is not uncommon to see requests for algorithms capable of handling problems containing  $10^5$ – $10^6$  examples (Smola & Schölkopf, 1998). The basic training algorithm (Vapnik, 1995) involves a solution to a quadratic programming problem, the size of which depends on the size of the training data. For the training data set of size  $l$ , the quadratic program for the pattern recognition SVM has  $l$  variables with non-negativity constraints,  $l$  inequality constraints and one equality constraint. The quadratic program for the regression SVM has  $2l$  variables with non-negativity constraints,  $2l$  inequality constraints and one equality constraint. In addition, the basic formulation of the SVM training algorithm requires storage of the kernel matrix of size  $l \times l$  or  $2l \times 2l$  respectively. Thus, the running time and storage of the kernel matrix are the two main optimization-related bottlenecks of SVM training algorithms.<sup>1</sup>

Original SVM implementations employed general-purpose optimization packages, such as MINOS, LOQO and others, that were not designed for problems of such size. It was soon discovered that these packages were not suitable for solutions to problems involving more

than a few hundred examples. The early special-purpose methods, proposed to speed-up training, have not brought much relief. Chunking (Vapnik, 1982) prescribes iteration through the training data, with accumulation of support vectors and addition of a “chunk” of new data until no more changes to the solution occur. The main problem with this method is that when the percentage of support vectors is high it essentially solves the problem of almost the full size more than once. Another method proposed in Kaufmann (1999) modifies traditional optimization algorithms (a combination of Newton’s and conjugate gradient methods) to yield the overall complexity of  $O(s^3)$  per iteration, where  $s$  is the (a priori unknown) number of support vectors. This can be a significant improvement over  $O(l^3)$ , however the number of support vectors is not guaranteed to be small.

Decomposition was the first practical method for solving large-scale SVM training problems. It was initially proposed for the pattern recognition SVM (Osuna, Freund, & Girosi, 1997) and subsequently extended to the regression SVM in Osuna (1998). The key idea of decomposition is to freeze all but a small number of optimization variables, and to solve a sequence of constant-size problems. The set of variables optimized at a current iteration is denoted as the *working set*. Because the working set is re-optimized, the value of the objective function is improved at each iteration, provided the working set is not optimal before re-optimization. Iteration is stopped when termination criteria, derived from Karush-Kuhn-Tucker (KKT) conditions, are satisfied to a required precision.

Selection of the working set is the most important issue in decomposition algorithms. First, the provision that the working set must be sub-optimal before re-optimization, is crucial to prevent the algorithm from cycling. Therefore, it is important to have criteria to test sub-optimality of any given working set. Second, working set selection affects the speed of the algorithm: if sub-optimal working sets are selected more or less at random the algorithm converges very slowly. Finally, working set selection is important in theoretical analysis of decomposition algorithms; in particular, in the recent convergence proof by Chang, Hsu, and Lin (1999).

The original decomposition algorithm essentially addressed only the first issue: the design of its termination criteria in all but pathological cases prevents composition of already optimal working sets. Its implementation featured some unpublished heuristics which provided reasonable convergence speed. Obviously, a formal framework for working set selection was highly desirable. Such framework can be provided by the classical method of feasible directions, proposed in the optimization theory in Zoutendijk (1960). The connection between this method and the working set selection problem was discovered by Joachims in a paper that has drawn wide attention (Joachims, 1999). However, Joachims’ algorithm is limited to the pattern recognition case because it uses the fact that the labels are  $\pm 1$ .

The main goal of the current article is to provide a unified treatment of the working set selection problem within the framework of the method of feasible directions. Specifically, we address, in a common way for the pattern recognition and the regression SVM, the following two issues:

- a criterion for identification of sub-optimal working sets,
- a heuristic for (approximately) optimal working set selection.

Resolution of these issues results in a new algorithm, termed the “maximal inconsistency algorithm”, which is applicable for both the pattern recognition and regression SVM.

Theoretical analysis of the new algorithm is aimed to justify the maximal inconsistency heuristic it uses. It is shown that for the pattern recognition SVM the new algorithm is equivalent to Joachims' algorithm, and for the regression SVM a similar algorithm exists based on the feasible direction principle. Similarly to the pattern recognition case, the maximal inconsistency algorithm satisfies theoretical requirements crucial for the proof of convergence. These relationships allow us to suggest the general notion of "feasible direction decomposition" algorithms.

Some experimental results are presented, intended to highlight certain properties of the maximal inconsistency algorithm:

- its scaling factors are consistent with the scaling factors of Joachims' algorithm, which corroborates the proved equivalence between the two.
- the algorithm has a linear convergence rate. This result is very important because it suggests the direction for theoretical convergence rate analysis, and also because it is consistent with the known linear convergence rate of gradient descent algorithms, to which the method of feasible directions belongs.
- profiled scaling factors demonstrate that complexity of the main steps in the algorithm is consistent with theoretical expectations. This pertains mostly to per-iteration complexity of certain computational steps.

It is beyond the scope of the article to present comprehensive numerical evaluation and comparison of feasible direction decomposition algorithms with other SVM training algorithms. The main competing algorithms are presented in the discussion, however numerical evaluation would require careful selection of benchmarks which are currently not established in the SVM community. Such evaluation, however, would be desirable in future.

The article is organized as follows. Section 2 provides basic formulations of decomposition algorithms for the pattern recognition and regression SVM. Section 3 presents the method of feasible directions. Sections 4 and 5 cover the issues related to feasible direction decomposition algorithms for the pattern recognition and regression SVM respectively. Exposition in these sections is carried out in different fashion. Section 4 is presented in a more or less chronological order. It starts with Joachims' algorithm, fully motivated by the feasible direction method. Then the criterion for testing optimality of a working set is presented, and the maximal inconsistency working set selection rule is derived from this criterion. Finally, equivalence between the maximal inconsistency algorithm and Joachims' algorithm is proved. The order of presentation is reversed for the regression SVM in Section 5. The maximal inconsistency algorithm is introduced first, followed by its interpretation as a feasible direction decomposition algorithm. Experimental results for the regression SVM are presented in Section 6.

## 2. SVM decomposition

Given the data  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_l, y_l)\}$ , and hence the kernel matrix  $K_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$ , the classical problem of SVM training can be expressed as the following quadratic

program:

$$\begin{aligned}
& \text{Maximize} && W(\tilde{\alpha}) = \tilde{\mathbf{e}}^T \tilde{\alpha} - \frac{1}{2} \tilde{\alpha}^T \tilde{H} \tilde{\alpha} \\
& \text{subject to:} && \tilde{\mathbf{c}}^T \tilde{\alpha} = 0 \\
& && \tilde{\alpha} - C\mathbf{1} \leq \mathbf{0} \\
& && \tilde{\alpha} \geq \mathbf{0}
\end{aligned} \tag{1}$$

The use of notation  $\tilde{x}$  for various quantities is aimed at highlighting the affinity between the two types of SVM. Interpretation of tilded quantities varies between the pattern recognition and the regression problems. For the pattern recognition SVM,

$$\tilde{\alpha} = \alpha \quad \tilde{\mathbf{e}} = \mathbf{1} \quad \tilde{H} = \mathbf{y}\mathbf{y}^T \odot K \quad \tilde{\mathbf{c}} = \mathbf{y} \tag{2}$$

For the regression SVM,

$$\tilde{\alpha} = \begin{bmatrix} \alpha \\ \alpha^* \end{bmatrix} \quad \tilde{\mathbf{e}} = \begin{bmatrix} \mathbf{y} - \epsilon\mathbf{1} \\ -\mathbf{y} - \epsilon\mathbf{1} \end{bmatrix} \quad \tilde{H} = \begin{bmatrix} K & -K \\ -K & K \end{bmatrix} \quad \tilde{\mathbf{c}} = \begin{bmatrix} \mathbf{1} \\ -\mathbf{1} \end{bmatrix} \tag{3}$$

Details of formulation of SVM training problems can be found in various sources, e.g., Vapnik (1995, 1999), Schölkopf (1997), and Smola (1998).

As it was mentioned in the introduction, the main idea of decomposition is to allow only a subset of optimization variables to change weights at a current iteration. The iteration process is repeated until termination conditions are met.

Let  $\tilde{\alpha}_B$  denote the variables included in the current working set (of fixed size  $q$ ), and let  $\tilde{\alpha}_N$  denote the rest of the variables. The corresponding parts of vectors  $\tilde{\mathbf{c}}$  and  $\tilde{\mathbf{y}}$  will also bear subscripts  $N$  and  $B$ . Matrix  $\tilde{H}$  will be partitioned into  $\tilde{H}_{BB}$ ,  $\tilde{H}_{BN} = \tilde{H}_{NB}^T$  and  $\tilde{H}_{NN}$ .

Optimization of the working set turns out to be also a quadratic program. This can be seen by rearranging the terms of the objective function and the equality constraint in (1) and dropping the terms independent of  $\tilde{\alpha}_B$  from the objective. The resulting quadratic program (the “sub-problem”) becomes:

$$\begin{aligned}
& \text{Maximize} && W_B(\tilde{\alpha}_B) = (\tilde{\mathbf{e}}_B^T - \tilde{\alpha}_N^T \tilde{H}_{NB}) \tilde{\alpha}_B - \frac{1}{2} \tilde{\alpha}_B^T \tilde{H}_{BB} \tilde{\alpha}_B \\
& \text{subject to:} && \tilde{\mathbf{c}}_B^T \tilde{\alpha}_B = -\tilde{\mathbf{c}}_N^T \tilde{\alpha}_N \\
& && \tilde{\alpha}_B - C\mathbf{1} \leq \mathbf{0} \\
& && \tilde{\alpha}_B \geq \mathbf{0}
\end{aligned} \tag{4}$$

The termination conditions of the decomposition algorithm are derived from the KKT conditions (Osuna, Freund, & Girosi, 1997; Osuna, 1998). It is best to consider them separately for the two types of SVM.

For the pattern recognition SVM, let<sup>2</sup>

$$Q_i = \sum_{j=1}^l \alpha_j y_j K(\mathbf{x}_i, \mathbf{x}_j) + b \tag{5}$$

where  $b$  is the threshold of SVM computed as (Joachims, 1999):

$$b = \frac{1}{|SV|} \sum_{i \in SV} \left[ y_i - \sum_{j=1}^l \alpha_j y_j K(\mathbf{x}_i, \mathbf{x}_j) \right] \quad (6)$$

Then the point satisfies the KKT conditions provided:

$$\begin{aligned} y_i q_i &= 1, & \text{if } 0 < \alpha_i < C \\ y_i q_i &\leq 1, & \text{if } \alpha_i = C \\ y_i q_i &\geq 1, & \text{if } \alpha_i = 0 \end{aligned} \quad (7)$$

For the regression SVM, let

$$q_i = \sum_{j=1}^l (\alpha_j - \alpha_j^*) K(\mathbf{x}_i, \mathbf{x}_j) + b \quad (8)$$

where

$$b = \frac{1}{|SV|} \sum_{i \in SV} \left[ y_i - \sum_{j=1}^l (\alpha_j - \alpha_j^*) K(\mathbf{x}_i, \mathbf{x}_j) - \epsilon \right] \quad (9)$$

Then the point satisfies the KKT conditions provided:

$$\begin{aligned} |q_i - y_i| &= \epsilon, & \text{if } 0 < \alpha_i < C \text{ and } 0 < \alpha_i^* < C \\ |q_i - y_i| &\geq \epsilon, & \text{if } \alpha_i = C \text{ or } \alpha_i^* = C \\ |q_i - y_i| &\leq \epsilon, & \text{if } \alpha_i = 0 \text{ and } \alpha_i^* = 0 \end{aligned} \quad (10)$$

Conditions (7) or (10) are checked for *individual points*: if, at a given iteration, one of them is violated for some point, this point is exchanged with some point in the current working set. Thus the new working set becomes sub-optimal, and strict improvement of the overall objective function is ensured.

The problem with using conditions (7) or (10) is that they require knowledge of the threshold  $b$ , which is also the Lagrange multiplier for the equality constraint of the SVM problem. The formulas for computing  $b$  given in (6) and (9) tacitly assume that the set  $SV$  of *unbounded* support vectors is non-empty. Usually this assumption is true; however, it cannot be guaranteed, especially for small working sets. While a simple trick can rectify the problem for the pattern recognition SVM, by letting

$$b = -\frac{m + M}{2} \quad (11)$$

where

$$m = \min_{i \in \text{SV}} \sum_{j=1}^l \alpha_j y_j K(\mathbf{x}_i, \mathbf{x}_j)$$

and

$$M = \max_{i \in \text{SV}} \sum_{j=1}^l \alpha_j y_j K(\mathbf{x}_i, \mathbf{x}_j)$$

a similar trick does not work for the regression SVM.

Overcoming this problem has been the initial motivation for the new termination condition proposed in this article. Instead of individual points, optimality of the entire working set will be considered. The new condition is described in detail in Section 4.2.

In general, intelligent selection of working sets is of primary importance for computational efficiency of decomposition algorithms. Before presenting in detail the working set selection techniques the source of their motivation needs to be introduced.

### 3. The method of feasible directions

Let  $\Omega$  be a feasible region of a general constrained optimization problem. Then a vector  $\mathbf{d}$  is said to be a *feasible direction* at the point  $\alpha \in \Omega$ , if there exists  $\bar{\lambda}$  such that  $\alpha + \lambda \mathbf{d} \in \Omega$  for all  $0 \leq \lambda \leq \bar{\lambda}$ .

The main idea of the method of feasible direction is to find a path from the initial feasible solution to the optimal solution by making steps along feasible directions. At each iteration the feasible direction algorithm proceeds as follows:

- find the optimal feasible direction—that is, the feasible direction providing the largest rate of increase of the objective function,
- determine the step length along the feasible direction that maximizes the objective function (“line search”).

The algorithm terminates when no feasible direction can be found which improves the objective function. For a general constrained optimization problem in the form

$$\begin{aligned} &\text{Maximize} && f(\alpha) \\ &\text{subject to:} && A\alpha \leq \mathbf{b} \end{aligned}$$

the optimal feasible direction is found by solving the direction finding linear program:<sup>3</sup>

$$\begin{aligned} &\text{Maximize} && \nabla f^T \mathbf{d} \\ &\text{subject to:} && A\mathbf{d} \leq 0 \\ &&& \|\mathbf{d}\|_2 \leq 1 \end{aligned} \tag{12}$$

The method of feasible directions can be applied directly to SVM training (Osuna, 1998). In this case, the respective optimal feasible direction problem can be stated as follows:

$$\text{Maximize } \tilde{\mathbf{g}}^T \mathbf{d} \quad (13)$$

$$\text{subject to: } \tilde{\mathbf{c}}^T \mathbf{d} = 0 \quad (14)$$

$$d_i \geq 0, \quad \text{if } \alpha_i = 0 \quad (15)$$

$$d_i \leq 0, \quad \text{if } \alpha_i = C \quad (16)$$

$$\|\mathbf{d}\|_2 \leq 1 \quad (17)$$

where  $\tilde{\mathbf{g}}$  is the gradient of the objective function of (1).

It turns out that solving the linear program of the full size at each iteration is expensive, and the overall performance of this method for SVM training is inferior to traditional optimization methods. However, with a slight modification, an approximate solution to the optimal feasible direction problem can be obtained in linear time, and it can serve as a powerful guide for working set selection. This approximate solution lies at the core of the feasible direction decomposition algorithms described in the ensuing sections.

#### 4. Feasible direction decomposition of the pattern recognition SVM

##### 4.1. Joachims' decomposition algorithm

The key observation of Joachims is that adding a requirement that only  $q$  components of  $\mathbf{d}$  be non-zero provides a straightforward working set selection rule: the variables corresponding to these non-zero components should be included in the new working set. Unlike Zoutendijk's method, the optimal feasible direction vector is not to be followed exactly. With re-optimization assumed to be cheap, one can afford finding an optimal solution *in the entire subspace* spanned by non-zero components of  $\mathbf{d}$ , instead of doing line search strictly along  $\mathbf{d}$ . Unfortunately, with the additional constraint, the optimal feasible direction problem becomes intractable. Therefore, one has to seek an approximate solution. One such solution is realized in Joachims' algorithm. In the rest of this section a detailed account of this solution is presented—in order to provide some insight and to underline the ideas used in the extension of this algorithm to the regression SVM.

The approximate solution is obtained by changing the normalization constraint (17) to

$$d_i \in \{-1, 0, 1\}$$

With this normalization, in order to satisfy the equality constraint (14),<sup>4</sup> it suffices that *the number of elements with sign matches between  $d_i$  and  $y_i$  is equal to the number of elements with sign mismatches between  $d_i$  and  $y_i$* . Obviously, for a working set of size  $q$  this condition holds only when each number is equal to  $q/2$ . Therefore, the equality constraint (14) can be enforced by performing two passes on the data: the "forward pass" selects  $q/2$  elements with sign mismatches, and the "backward pass"— $q/2$  elements with sign matches.<sup>5</sup>

How should directions be determined for the elements? Recall that the goal is to maximize the objective function (13) subject to constraints. In the pattern recognition case the gradient

of the objective is given by:

$$\mathbf{g} = \mathbf{1} - K\boldsymbol{\alpha} \quad (18)$$

In the absence of constraints (14)–(16) the maximum of the objective function would have been achieved by selecting  $q$  points with the highest values of  $|g_i|$  and assigning directions  $d_i = -\text{sign}(g_i)$  to them. Consider the largest contribution  $\gamma_k$  to the objective function provided by some point  $k$ , *subject to the equality constraint*. During the forward pass the signs of  $d_k$  and  $y_k$  must be different, therefore:

$$\begin{aligned} y_k = 1 &\Rightarrow d_k = -1 \Rightarrow \gamma_k = \max_{i: y_i=1} (-g_i) \Rightarrow \gamma_k = \min_{i: y_i=1} (g_i y_i) \\ y_k = -1 &\Rightarrow d_k = 1 \Rightarrow \gamma_k = \max_{i: y_i=-1} (g_i) \Rightarrow \gamma_k = \min_{i: y_i=-1} (g_i y_i) \end{aligned}$$

and hence, combining the subscripts,

$$\gamma_k = \min_i (g_i y_i)$$

Likewise, for the backward pass the signs of  $d_k$  and  $y_k$  must be the same, therefore:

$$\begin{aligned} y_k = 1 &\Rightarrow d_k = 1 \Rightarrow \gamma_k = \max_{i: y_i=1} (g_i) \Rightarrow \gamma_k = \max_{i: y_i=1} (g_i y_i) \\ y_k = -1 &\Rightarrow d_k = -1 \Rightarrow \gamma_k = \max_{i: y_i=-1} (-g_i) \Rightarrow \gamma_k = \max_{i: y_i=-1} (g_i y_i) \end{aligned}$$

and hence, combining the subscripts,

$$\gamma_k = \max_i (g_i y_i)$$

Thus the quantity  $g_i y_i$  reflects the element's contribution to the objective function subject to the equality constraint. The working set composition rule can then be stated as follows: sort data elements by  $g_i y_i$  in increasing order and select  $q/2$  elements from the front of the list (hence, the "forward pass"), and  $q/2$  elements from the back of the list (hence, the "backward pass").<sup>6</sup>

Finally, to account for inequality constraints (15) and (16), some points may have to be skipped if they violate one of these constraints. For such points, the direction leading to improvement of the objective function of the optimal feasible direction problem is infeasible.

Joachims' algorithm is summarized in Algorithm 1.

---

**Algorithm 1** Joachims' SVM decomposition algorithm.

---

Let  $S$  be a list of all samples.

**while** (termination conditions (7) are not met)

- sort  $S$  by  $y_i g_i$  in increasing order
  - select  $q/2$  samples from the front of  $S$  such that  $0 < \alpha_i < C$  or  $d_i = -y_i$  satisfies (15) and (16) (forward pass)
  - select  $q/2$  samples from the back of  $S$  such that  $0 < \alpha_i < C$  or  $d_i = y_i$  satisfies (15) and (16) (backward pass)
  - re-optimize the working set
-



#### 4.2. Optimality of a working set

As it was mentioned earlier, the point-wise termination criteria in Osuna's and Joachims' algorithms require knowledge of the threshold  $b$  of the SVM. This threshold can be difficult to calculate, especially for the regression SVM. In this section an alternative termination condition is presented, which allows testing whether or not the entire working set is sub-optimal and hence suitable for re-optimization. The new condition is based on examination of the KKT conditions. We will first state the KKT theorem in the general form for the problem (1) and will specialize it later for the particular kind of SVM.

**Theorem 1** (*Karush-Kuhn-Tucker Theorem*). *The primal vector  $\tilde{\alpha}$  solves the quadratic problem (1) if and only if it satisfies the constraints of (1) and there exists a dual vector  $\mathbf{u}^T = (\mu \mathbf{\Pi}^T \mathbf{\Upsilon}^T)$  such that:*

$$\tilde{H} \tilde{\alpha} - \tilde{\mathbf{e}} + \mu \tilde{\mathbf{c}} + \mathbf{\Upsilon} - \mathbf{\Pi} = \mathbf{0} \quad (19)$$

$$\mathbf{\Upsilon} \geq \mathbf{0} \quad (20)$$

$$\mathbf{\Pi} \geq \mathbf{0} \quad (21)$$

$$\mathbf{\Upsilon}^T (\tilde{\alpha} - C \mathbf{1}) = 0 \quad (22)$$

$$\mathbf{\Pi}^T \tilde{\alpha} = 0 \quad (23)$$

Conditions (19) and (21) can be combined into an equivalent condition:

$$\mathbf{\Pi} = \tilde{H} \tilde{\alpha} - \tilde{\mathbf{e}} + \mu \tilde{\mathbf{c}} + \mathbf{\Upsilon} \geq \mathbf{0} \quad (24)$$

It follows from the Karush-Kuhn-Tucker Theorem that if for all  $\mathbf{u}$  satisfying conditions (20)–(23) the system of inequalities (24) is inconsistent then the solution of problem (1) is *not* optimal. Since the sub-problem (4) was obtained by merely re-arranging terms in the objective function and the constraints of the initial problem (1), the same conditions guarantee that the sub-problem (4) is not optimal. Thus, our main strategy for identifying sub-optimal working sets will be to enforce inconsistency of the system (24) while satisfying conditions (20)–(23). Notice that all constant terms in (24) represent the negative gradient vector  $-\tilde{\mathbf{g}} = \tilde{H} \tilde{\alpha} - \tilde{\mathbf{e}}$ .

Let us now consider the particular form of the pattern recognition SVM. Each inequality in (24) can be written as follows:

$$\pi_i = -g_i + y_i \mu + v_i \geq 0 \quad (25)$$

where  $g_i$  is defined as in (18). Consider three cases, according to the values that  $\alpha_i$  can take:

1.  $0 < \alpha_i < C$ . In this case, the complementarity conditions (22), (23) imply that  $\pi_i = 0$ ,  $v_i = 0$ . Then inequality (25) becomes

$$-g_i + y_i \mu = 0 \quad (26)$$

2.  $\alpha_i = 0$ . In this case, the complementarity condition (22) implies that  $\nu_i = 0$ . Then inequality (25) becomes

$$-g_i + y_i \mu \geq 0 \quad (27)$$

3.  $\alpha_i = C$ . In this case, the complementarity condition (23) implies that  $\pi_i = 0$ . Then inequality (25) becomes

$$-g_i + y_i \mu \leq 0 \quad (28)$$

It can be easily seen that enforcing the complementarity constraints (22), (23) causes  $\mu$  to become the only free variable in the system (24). Each point restricts  $\mu$  to a certain interval on a real line. Such intervals will be denoted as  $\mu$ -sets in the rest of the article. The rule for computation of  $\mu$ -sets can be summarized as follows:

$$\mathcal{M}_i = \begin{cases} [g_i/y_i], & \text{if } 0 < \alpha_i < C \\ [g_i, +\infty), & \text{if } \alpha_i = 0 \text{ and } y_i = 1 \\ (-\infty, -g_i], & \text{if } \alpha_i = 0 \text{ and } y_i = -1 \\ (-\infty, g_i], & \text{if } \alpha_i = C \text{ and } y_i = 1 \\ [-g_i, +\infty), & \text{if } \alpha_i = C \text{ and } y_i = -1 \end{cases} \quad (29)$$

The development in this section proves the following theorem:

**Theorem 2.** *The vector  $\alpha$  solves the quadratic program (1)–(2) if and only if intersection of  $\mu$ -sets computed by (29) is non-empty.*

It also follows from the expressions (29) that if at least one  $\alpha_i$  is strictly between the bounds, then at the optimal solution the intersection of all  $\mu$ -sets (which is non-empty by the Kuhn-Tucker theorem) is a single point set. This is consistent with the known property of SVM that at the optimal solution  $b = \mu$ . The intersection of  $\mu$ -sets at the optimal solution can only be a non-empty, non single point set if all variables are at the bounds. In this case any point from the intersection of  $\mu$ -sets can be taken as  $b$ ; in particular, it can be the value suggested in (11).

#### 4.3. Maximal inconsistency algorithm

While inconsistency of the working set at each iteration guarantees convergence of decomposition, the rate of convergence is quite slow if arbitrary inconsistent working sets are chosen. A natural heuristic is to select “maximally inconsistent” working sets, hoping that such choice would provide the greatest improvement of the objective function. The notion of “maximal inconsistency” is easy to define: let it be the gap between the smallest right boundary and the largest left boundary of  $\mu$ -sets of elements in the training set:

$$G = L - R \\ L = \max_{0 < i < l} \mu_i^l, \quad R = \min_{0 < i < l} \mu_i^r$$

where  $\mu_i^l, \mu_i^r$  are the left and the right boundaries respectively (possibly minus or plus infinity) of  $\mu$ -set  $\mathcal{M}_i$ . It is convenient to require that the largest possible inconsistency gap be maintained between all pairs of points comprising the working set. The obvious implementation of this strategy is to select  $q/2$  elements with the largest values of  $\mu^l$  and  $q/2$  elements with the smallest values of  $\mu^r$ .

One feature of Joachims' method needs to be retained—rejection of Zoutendijk-infeasible points (c.f. (15), (16)). Inclusion of such points into the working set doesn't make sense anyway because their values of  $\alpha$  will not change after re-optimization. In the pattern recognition case, the  $\mu$ -set constraints are not capable of encoding feasibility.<sup>7</sup> Since the notion of the direction is not explicitly maintained in the maximal inconsistency algorithm the feasibility test needs to be slightly modified: the point is infeasible if  $\alpha_i = 0$  and  $g_i < 0$  or if  $\alpha_i = C$  and  $g_i < 0$ .

The maximal inconsistency strategy is summarized in Algorithm 2.

#### 4.4. Equivalence between Joachims' algorithm and the maximal inconsistency algorithm

So far the motivation for the maximal inconsistency algorithm has been purely heuristic. A natural question arises: does inconsistency gap indeed provide a good measure of optimality for a working set at a given iteration? An affirmative answer to it is developed in this section by showing the equivalence between Joachims' algorithm and the maximal inconsistency algorithm.

---

**Algorithm 2** Maximal inconsistency algorithm for the pattern recognition SVM.

---

Let  $S$  be a list of all samples.

**while** (intersection of all  $\mu$ -sets is empty)

- compute  $\mathcal{M}_i$  according to the rules (29) for all elements in  $S$
  - select  $q/2$  feasible samples with the largest values of  $\mu^l$  (“left pass”)
  - select  $q/2$  feasible samples with the smallest values of  $\mu^r$  (“right pass”)
  - re-optimize the working set
- 

To show that the algorithms are equivalent we have to prove that they produce identical working sets at each iteration<sup>8</sup> and that their termination conditions are equivalent. The two propositions below handle each claim.

**Proposition 1.** *The working set of Joachims' algorithm is identical to the working set of the maximal inconsistency algorithm at each iteration.*

**Proof:** The statement above will be proved for the half of the working set; namely, that the set of elements selected by the “forward pass” of Joachims' algorithm is identical to the set of elements selected by the “right pass” of the maximal inconsistency algorithm. Similar argument allows to establish equivalence of the other half of working sets.

Let  $F$  be the set of all feasible samples at some iteration. Let  $r_J(p)$  denote the rank (position) of sample  $p$  in the array obtained by sorting  $F$  by  $y_i g_i$ ; let  $r_I(p)$  be the rank of sample  $p$  in the array obtained by sorting  $F$  by  $\mu_i^r$ .

Let the set  $H_J = \{p : r_J(p) \leq q/2\}$  be the half of the working set selected by the “forward pass” of Joachims’ algorithm. Let  $\bar{p}$  be the sample whose  $r_J(\bar{p}) = q/2$ , i.e. the element in  $H_J$  with the largest value of the key. Let  $H_I = \{p : \mu^r(p) \leq \mu^r(\bar{p})\}$ . We are going to prove that  $H_J \equiv H_I$ .

For any sample  $p_i$  selected by the “forward pass”  $d_i = -y_i$ . Considering the possible values of  $\alpha_i$ , we conclude that, for any  $p \in H_J$ :

$$\begin{aligned} 0 < \alpha_i < C & \Rightarrow \mu_i^r = y_i g_i \\ \alpha_i = 0 & \Rightarrow d = 1 \Rightarrow y_i = -1 \Rightarrow \mu_i^r = -g_i = -1 \cdot g_i = y_i g_i \\ \alpha_i = C & \Rightarrow d = 1 \Rightarrow y_i = 1 \Rightarrow \mu_i^r = g_i = 1 \cdot g_i = y_i g_i \end{aligned} \quad (30)$$

We can see that the order of  $\mu^r(p)$  is preserved when mapping  $H_J \mapsto H_I$ . Therefore, if  $p \in H_J$  then  $p \in H_I$ .

To prove set equivalence, it remains to be shown that if  $p \in H_I$  then  $p \in H_J$ . Suppose, by the way of contradiction this is not the case, i.e., there exists a sample  $\tilde{p}$  such that  $\mu^r(\tilde{p}) < \mu^r(\bar{p})$  and  $\tilde{p} \notin H_J$ . Since  $\tilde{p} \notin H_J$ ,  $r_J(\tilde{p}) \geq r_J(\bar{p})$ . What is  $\mu^r(\tilde{p})$  equal to? Three of the five possible cases are covered in (30), and in each of them  $\mu^r(\tilde{p}) = g_i y_i$ . In this case  $\mu^r(\tilde{p}) < \mu^r(\bar{p}) \Rightarrow r_J(\tilde{p}) < r_J(\bar{p})$  which contradicts our previous conclusion that  $r_J(\tilde{p}) \geq r_J(\bar{p})$ . In the remaining two cases of (29)  $\mu^r(\tilde{p}) = +\infty$  which contradicts our assumption that  $\mu^r(\tilde{p}) < \mu^r(\bar{p})$ . Thus we can conclude that  $H_J \equiv H_I$ .  $\square$

**Proposition 2.** *Termination conditions of the Joachims’ algorithm are satisfied if and only if the termination conditions of the maximal inconsistency algorithm are satisfied.*

**Proof:** The maximal inconsistency algorithm terminates when the system (24) is consistent while at the same time conditions (20)–(23) are enforced. Hence, the KKT conditions are satisfied, and consequently, the algorithm terminates if and only if the optimal solution is found. Likewise, termination conditions (5)–(7) also have an “if and only if” relationship to the KKT conditions and hence to the optimality of the solution, except when the solution contains all variables strictly at the bounds and (11) is used for calculation of  $b$ . In the latter case, however, condition (11) satisfies the KKT conditions of the *primal* SVM training problem (to which the problem (1)–(2) is a dual). Then it follows from Dorn’s duality theorem (Mangasarian, 1969), that the solution of the dual problem is also optimal. Hence, both algorithms terminate at the same point in the solution space.  $\square$

## 5. Feasible direction decomposition of the regression SVM

### 5.1. Maximal inconsistency algorithm

We now turn our attention to the maximal inconsistency algorithm for the regression SVM. Recall that the quadratic program for the latter is given by Eqs. (1) and (3). The derivation

will progress in the same way as in the pattern recognition case and will consist of: (a) the statement of the Karush-Kuhn-Tucker Theorem, (b) the derivation of rules for computation of  $\mu$ -sets, and (c) the definition of the inconsistency gap to be used for working set selection in the algorithm.

The general statement of the Karush-Kuhn-Tucker Theorem, as well as its use to test optimality of the working set, remain the same (see Theorem 1 and the ensuing discussion).

For the regression SVM, each inequality in (24) has one of the following forms:

$$\pi_i = -\phi_i + \epsilon + v_i + \mu \geq 0 \quad (31)$$

$$\pi_i^* = \phi_i + \epsilon - v_i^* - \mu \geq 0 \quad (32)$$

where

$$\phi_i = y_i - \sum_{j=1}^l (\alpha_j - \alpha_j^*) K_{ij} \quad (33)$$

Considering the possible values of  $\alpha_i$  we have:

1.  $\alpha_i = 0$ . In this case, by complementarity condition (22)  $v_i = 0$ . Then inequality (31) becomes:

$$\pi_i = -\phi_i + \epsilon + \mu \geq 0 \Rightarrow \mu \geq \phi_i - \epsilon$$

2.  $\alpha_i = C$ . By complementarity condition (23)  $\pi_i = 0$ . Then inequality (31) becomes:

$$-\phi_i + \epsilon + \mu + v_i = 0 \Rightarrow \mu \leq \phi_i - \epsilon$$

3.  $0 < \alpha_i < C$ . By complementarity conditions (22) (23),  $v_i = 0$ ,  $\pi_i = 0$ . Then inequality (31) becomes:

$$-\phi_i + \epsilon + \mu = 0 \Rightarrow \mu = \phi_i - \epsilon$$

Similar reasoning for  $\alpha_i^*$  and inequality (32) yields the following results:

1.  $\alpha_i^* = 0$ . Then

$$\mu \leq \phi_i + \epsilon$$

2.  $\alpha_i^* = C$ . Then

$$\mu \geq \phi_i + \epsilon$$

3.  $0 < \alpha_i^* < C$ . Then

$$\mu = \phi_i + \epsilon$$

Finally, taking into account that for the regression SVM  $\alpha_i \alpha_i^* = 0$ , the rules for computation of  $\mu$ -sets for the regression SVM are the following:

$$\mathcal{M}_i = \begin{cases} [\phi_i - \epsilon, \phi_i + \epsilon], & \text{if } \alpha_i = 0, \alpha_i^* = 0 \\ [\phi_i - \epsilon, \phi_i - \epsilon], & \text{if } 0 < \alpha_i < C, \alpha_i^* = 0 \\ (-\infty, \phi_i - \epsilon], & \text{if } \alpha_i = C, \alpha_i^* = 0 \\ [\phi_i + \epsilon, \phi_i + \epsilon], & \text{if } \alpha_i = 0, 0 < \alpha_i^* < C \\ [\phi_i + \epsilon, +\infty), & \text{if } \alpha_i = 0, \alpha_i^* = C \end{cases} \quad (34)$$

This development leads to the following theorem.

**Theorem 3.** *The vector  $\alpha$  solves the quadratic program (1)–(3) if and only if intersection of  $\mu$ -sets computed by (34) is non-empty.*

The maximal inconsistency algorithm for the regression SVM is summarized in Algorithm 3.

Feasibility of samples can be tested by the following rule: during the left pass skip samples with  $\alpha_i = C, \alpha_i^* = 0$ ; during the right pass skip samples with  $\alpha_i = 0, \alpha_i^* = C$ . Justification of this rule is given in Lemma 1 in Section 5.2.

### 5.2. Interpretation of the maximal inconsistency algorithm in the feasible direction framework

As it was shown in Section 4.4, the maximal inconsistency algorithm is equivalent to Joachims' algorithm, which was motivated by Zoutendijk's feasible direction problem. In this section it will be shown that the maximal inconsistency algorithm for the regression SVM can also be interpreted as a feasible direction algorithm.

---

**Algorithm 3** Maximal inconsistency algorithm for the regression SVM.

---

Let  $S$  be a list of all samples.

**while** (intersection of  $\mu$ -sets is empty)

- compute  $\mathcal{M}_i$  according to the rules (34) for all elements in  $S$
  - select  $q/2$  feasible samples with the largest values of  $\mu^l$  (“left pass”)
  - select  $q/2$  feasible samples with the smallest values of  $\mu^r$  (“right pass”)
  - re-optimize the working set
- 

Recall the SVM optimal feasible direction problem stated in (13)–(17). The problem-specific components  $\tilde{\mathbf{g}}$  and  $\tilde{\mathbf{c}}$  have the following expressions for the regression SVM:

$$\tilde{\mathbf{g}} = \begin{bmatrix} \mathbf{g} \\ \mathbf{g}^* \end{bmatrix}, \quad \tilde{\mathbf{c}} = \begin{bmatrix} \mathbf{1} \\ -\mathbf{1} \end{bmatrix} \quad (35)$$

where

$$\begin{aligned} g_i &= \phi_i - \epsilon \\ g_i^* &= -\phi_i - \epsilon \end{aligned} \quad (36)$$

and  $\phi$  is defined in (33). In addition, a feasible direction algorithm must satisfy the constraint

$$|\{d_i : d_i \neq 0\}| = q \quad (37)$$

To develop an equivalent feasible direction algorithm, we will construct a mapping  $\Phi_{\bar{\alpha}}$ , which maps the state of the maximal inconsistency algorithm to a direction vector  $\mathbf{d}$ , and a normalization  $N_{\mathbf{d}}$ , similar to Joachims' normalization, which replaces (17). The construction will possess the following properties:

1. At each iteration,  $\mathbf{d} = \Phi_{\bar{\alpha}}$  is the solution to the optimal feasible direction problem with normalization  $N_{\mathbf{d}}$ .
2. Termination condition of the maximal inconsistency algorithm holds if and only if  $\mathbf{d}$  is a zero direction.

Intuitively, the first property shows that working sets selected by the maximal inconsistency algorithm are the same as those selected by a feasible direction algorithm using normalization  $N_{\mathbf{d}}$ . The second property ensures that both algorithms terminate at the same time.

Consider the normalization

$$d_i \in \{(1, 0), (0, -1), (-1, 0), (0, 1), (0, 0)\} \quad (38)$$

and the mapping

$$(\Phi_{\bar{\alpha}})_i = \begin{cases} (1, 0) \text{ or } (0, -1), & \text{whichever is feasible, during the left pass} \\ (-1, 0) \text{ or } (0, 1), & \text{whichever is feasible, during the right pass} \\ (0, 0), & \text{if sample } i \text{ is infeasible or not reached} \end{cases} \quad (39)$$

For the sake of brevity, the optimal feasible direction problem comprising Eqs. (13)–(16), (37), (38) will be denoted as the *feasible direction problem*.

First, we need to make sure that  $\Phi_{\bar{\alpha}}$  is not ambiguous, i.e., that only one of the non-zero directions suggested in (39) is feasible.

**Lemma 1.** *Mapping  $\Phi_{\bar{\alpha}}$  is not ambiguous.*

**Proof:** Let us denote directions  $(1, 0)$ ,  $(0, -1)$ ,  $(-1, 0)$ ,  $(0, 1)$  as Ia, Ib, IIa, IIb (Type I directions are assigned during the left pass, and Type II directions—during the right pass). Table 1 shows feasibility of different directions depending on the values of optimization variables: Infeasibility of directions marked by  $\dagger$  is due to the special property of the regression SVM that  $\alpha_i \alpha_i^* = 0$ . It follows from the table that at each pass only one direction is feasible.  $\square$

Table 1. Feasibility of directions.

Optimization variables	Feasible	Infeasible
$\alpha_i = 0, \alpha_i^* = 0$	Ia, IIb	Ib, IIa
$0 < \alpha_i < C, \alpha_i^* = 0$	Ia, IIa	Ib, IIb <sup>†</sup>
$\alpha_i = C, \alpha_i^* = 0$	IIa	Ia, Ib, IIb <sup>†</sup>
$\alpha_i = 0, 0 < \alpha_i^* < C$	Ib, IIb	Ia <sup>†</sup> , IIa
$\alpha_i = 0, \alpha_i^* = C$	Ib	Ia <sup>†</sup> , IIa, IIb

This lemma justifies the feasibility test of the maximal inconsistency algorithm. It is clear from Table 1 that during the left pass there is no feasible direction for samples with  $\alpha_i = C, \alpha_i^* = 0$ . Likewise, during the right pass there is no feasible direction for samples with  $\alpha_i = 0, \alpha_i^* = C$ .

The next two lemmas show that  $\Phi_{\tilde{\alpha}}$  provides a solution to the feasible direction problem.

**Lemma 2.**  $\Phi_{\tilde{\alpha}}$  satisfies all constraints of the feasible direction problem.

**Proof:** The equality constraint of the optimal feasible direction problem for the regression SVM has the form:

$$\sum_{i=1}^l d_i = \sum_{i=1}^l d_i^* \quad (40)$$

Let  $m_1, m_2, n_1, n_2$  be the number of selected elements with directions  $\{1, 0\}, \{-1, 0\}, \{0, -1\}, \{0, 1\}$  respectively. Then

$$\begin{aligned} \sum_{i=1}^l d_i &= m_1 \cdot 1 + m_2 \cdot (-1) = m_1 - m_2 \\ \sum_{i=1}^l d_i^* &= n_1 \cdot (-1) + n_2 \cdot 1 = n_2 - n_1 \end{aligned}$$

By the selection policy:

$$m_1 + n_1 = q/2 = m_2 + n_2$$

from which it follows that

$$m_1 - m_2 = n_2 - n_1$$



Hence,  $\Phi_{\bar{\alpha}}$  satisfies the equality constraint of the optimal feasible direction problem.

Inequality constraints and the cardinality constraint are trivially satisfied by the construction of  $\Phi_{\bar{\alpha}}$ .  $\square$

**Lemma 3.**  $\Phi_{\bar{\alpha}}$  provides the optimal value of the objective function of the feasible direction problem.

**Proof:** Let  $B_l$  and  $B_r$  denote the halves of the working set selected by the left and the right passes respectively. Suppose, by the way of contradiction, that there exists a feasible sample  $i$  such that  $g_i d_i > g_j d_j$  for some sample  $j \in B_l \cup B_r$ .

For any element  $k$  considered during the left pass

$$g_k d_k = \begin{cases} 1 \cdot g_k = \phi_k - \epsilon = \mu_k^l & \text{if } (1, 0) \text{ is feasible} \\ -1 \cdot g_k^* = -(-\phi_k - \epsilon) = \phi_k + \epsilon = \mu_k^l & \text{if } (0, -1) \text{ is feasible} \end{cases}$$

Therefore, if  $j \in B_l$ , then

$$g_i d_i = \mu_i^l < \mu_j^l = g_j d_j$$

which contradicts the hypothesis that  $g_i d_i > g_j d_j$ .

Likewise, for any element  $k$  considered during the right pass

$$g_k d_k = \begin{cases} -1 \cdot g_k = -(\phi_k - \epsilon) = -\mu_k^r & \text{if } (-1, 0) \text{ is feasible} \\ 1 \cdot g_k^* = -\phi_k - \epsilon = -(\phi_k + \epsilon) = -\mu_k^r & \text{if } (0, 1) \text{ is feasible} \end{cases}$$

Therefore, if  $j \in B_r$ , then

$$-g_i d_i = \mu_i^r > \mu_j^r = -g_j d_j$$

which contradicts the hypothesis that  $g_i d_i > g_j d_j$ .  $\square$

**Lemma 4.** Intersection of  $\mu$ -sets is non-empty if and only if the feasible direction problem has a zero solution.

**Proof:** By Theorem 3, non-empty intersection of  $\mu$ -sets implies optimality of the solution to the quadratic program (1)–(3) which rules out existence of a non-zero feasible direction which would otherwise have led to a new optimal solution.

On the other hand, if the optimal solution to the feasible direction problem is zero, this implies that all other feasible directions have negative projections on the gradient vector, and hence decrease the value of the objective function. It follows that the solution of the quadratic program (1)–(3) is optimal, and intersection of  $\mu$ -sets is non-empty.  $\square$

Lemmas 1–4 prove the two properties of the mapping  $\Phi_{\bar{\alpha}}$  and the normalization  $N_{\mathbf{d}}$  claimed earlier in this section.

## 6. Experimental results

The aim of this section is to provide insight into some properties of feasible direction decomposition algorithms that might explain their behavior in different situations. It is not our intent to perform comprehensive numerical evaluation of several SVM training algorithms. The following issues are of prime interest:

- Overall scaling factors. This is the traditional way of analyzing performance of SVM training algorithms, introduced by Platt (1999) and Joachims (1999). To perform at least qualitative comparison with their results (which are reported for the pattern recognition SVM, while we are primarily interested in the regression case), a similar evaluation is performed for the maximal inconsistency algorithm.
- Experimental convergence rates. In numerical optimization literature algorithms are most often evaluated in terms of their convergence rates. Since decomposition algorithms are iterative and derive their core ideas from optimization theory, it is natural to attempt to establish their rates of convergence. It will be shown that the maximal inconsistency algorithm seems to have linear rate of convergence, which is consistent with known linear convergence rates for gradient descent methods. This result is particularly important in the context of ongoing theoretical investigation of convergence of decomposition algorithm.
- Profiled scaling factors. While decreasing the number of iterations is highly desirable, it must not be achieved at a cost of significantly increasing the cost of an iteration. It is therefore important to investigate the “profile” of one iteration of the decomposition algorithm. The results will demonstrate that complexity of the main steps of the algorithm is linear, much as predicted theoretically.
- Optimal working set size. Since the working set size is an additional parameter of the feasible direction algorithms, it is desirable to investigate its impact on the training time. It turns out that dependence of training time on the working set size is tightly related to iteration profiles.

Experimental evaluation of the new algorithm was performed on the modified KDD Cup 1998 data set. The original data set is available under <http://www.ics.uci.edu/~kdd/databases/kddcup98/kddcup98.html>.

The following modifications were made to obtain a pure regression problem:

- All 75 character fields were eliminated.
- Numeric fields CONTROLN, ODATEDW, TCODE and DOB were eliminated.

The remaining 400 features and the labels were scaled between 0 and 1. Initial subsets of the training database of different sizes were selected for evaluation of the scaling properties of the new algorithm. Experiments were run on a SUN4U/400 Ultra-450 workstation with 300 MHz clock and 2048M RAM. RBF kernel with  $\gamma = 10$ , termination accuracy<sup>9</sup> 0.001, and cache size of 300 M were used. The value of the box constraint  $C$  was 1, and the working set of size 10 was used. Two sets of experiments were performed: one using the full set of 400 features, another one using only the first 50 features (“reduced set”). As it turns out, the

second problem is more constrained, with a larger proportion of bounded support vectors. Also, for the full set of features kernel computation dominates the overall training time.

The current implementation does not employ the shrinking heuristic described in Joachims (1999). It uses MINOS to solve optimization (sub)-problems.

### 6.1. Overall scaling factors

The training times, with and without decomposition,<sup>10</sup> for different samples sizes, are displayed in Tables 2 and 3, for the full and the reduced sets of features respectively. The scaling factors are computed by plotting training times versus sample sizes on a log-log scale and fitting straight lines. The SV-scaling factors are obtained in the same fashion, only using the number of unbounded support vectors instead of the sample size as an abscissa. The actual plots are shown in figure 1. The plots demonstrate fine linear dependence on a log-log scale.

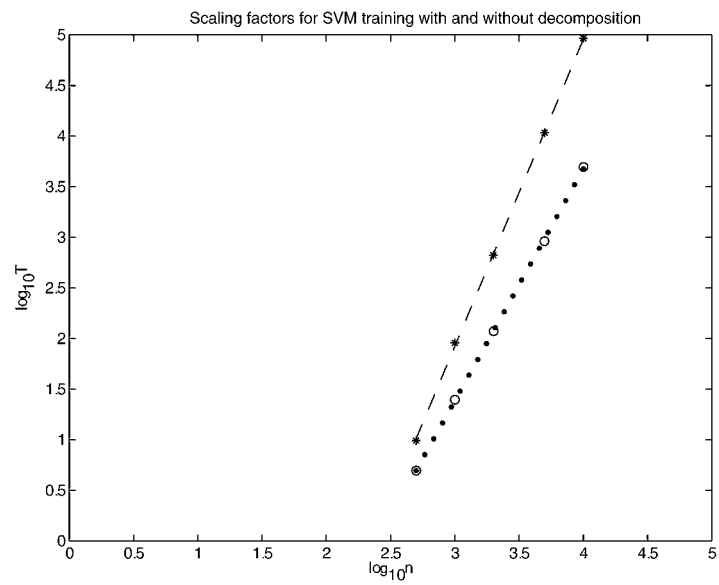
It can be easily seen that decomposition improves the running time by an order of magnitude. Its scaling factors are also significantly better. Although this results is far from

Table 2. Training time (sec) and number of SVs for the KDD Cup problem.

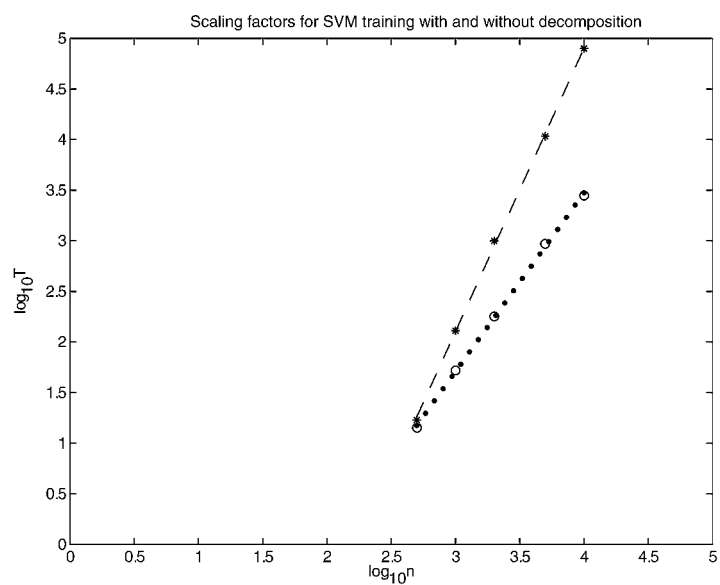
Examples	No dcmp			dcmp		
	Time	Total SV	BSV	Time	Total SV	BSV
500	9.818	184	0	4.962	170	0
1000	91.193	454	0	24.907	429	0
2000	665.566	932	2	118.131	894	2
5000	10785.4	2305	7	914.359	2213	7
10000	92684.7	4598	28	4958.26	4454	26
Scaling factor			3.03			2.29
SV-scaling factor			2.86			2.13

Table 3. Training time (sec) and number of SVs for the KDD Cup problem, reduced feature set.

Examples	No dcmp			dcmp		
	Time	Total SV	BSV	Time	Total SV	BSV
500	16.958	114	29	14.175	128	26
1000	129.591	242	65	52.451	324	59
2000	998.191	445	114	178.752	656	104
5000	10759.6	977	323	929.523	1667	255
10000	79323.2	1750	633	2782.35	3383	491
Scaling factor			2.80			1.76
SV-scaling factor			3.12			1.60



(a)



(b)

Figure 1. Scaling factor fits: (a) full set of features (Table 2), (b) reduced set of features (Table 3).

Table 4. Objective function values for KDD Cup problem.

Examples	No dcmp	dcmp	Ratio
500	0.10799	0.10795	0.99962
1000	0.43501	0.43488	0.99970
2000	1.27695	1.27662	0.99974
5000	3.19561	3.19414	0.99953
10000	6.74862	6.74507	0.99947

Table 5. Objective function values for KDD Cup problem, reduced feature set.

Examples	No dcmp	dcmp	Ratio
500	1.11129	1.10823	0.99724
1000	3.25665	3.24853	0.99750
2000	6.41472	6.39787	0.99456
5000	14.5490	14.5164	0.99776
10000	26.9715	26.9119	0.99779

surprising, noteworthy is consistency of the observed scaling factors with the scaling factors presented in Joachims (1999) for the pattern recognition SVM. This supports the claim of equivalence between the maximal inconsistency and Joachims' algorithms made earlier in this article.

A number of other interesting findings can be made from the results above. First, training with decomposition does not produce an identical solution to training without. The solutions differ in the number of support vectors, especially for the more constrained problem with the reduced set of features. This difference is due to the fact that termination conditions of the decomposition algorithm require that the KKT condition is satisfied only to a given numerical precision. Thus the decomposition algorithm does have a disadvantage of producing an approximate solution, which can be further seen from Tables 4 and 5. These tables display the values of the objective functions attained by training with and without decomposition, and their ratio. The latter shows that in the relative terms a solution of the more constrained problem with the reduced feature set is roughly 10 times worse than the solution of the less constrained problem with the full feature set. However, no deterioration of the accuracy with the sample size has been observed.

Another important observation is that the scaling factors and the SV-scaling factors vary among different problems. For the two particular problems, a possible explanation might be that the fixed termination accuracy is in fact looser for the more constrained problem, thereby producing a less accurate solution but taking less time in doing so. In general, however, the results above demonstrate that the scaling factors produce a rather crude measure of performance of the decomposition algorithms.<sup>11</sup>

## 6.2. Convergence rates

In optimization literature a common performance measure for iterative algorithms is convergence rate. The notion of convergence rate is generally defined for any numeric sequence; for the purpose of analysis of decomposition algorithms we will be concerned with the sequence of objective function values. The following definitions can be found in Nocedal and Wright (1999).

Let  $x_k$  be a sequence in  $\mathbb{R}^n$  that converges to  $x^*$ . Convergence is said to be *Q-linear* if there is a constant  $r \in (0, 1)$  such that

$$\frac{\|x_{k+1} - x^*\|}{\|x_k - x^*\|} \leq r, \quad \text{for all } k \text{ sufficiently large.}$$

The prefix “Q” stands for “quotient” because the quotient of successive distances from the limit point is considered. Likewise, convergence is said to be *Q-super-linear* if

$$\lim_{k \rightarrow \infty} \frac{\|x_{k+1} - x^*\|}{\|x_k - x^*\|} = 0,$$

and is said to be of *order*  $p$  (for  $p > 1$ ) or *quadratic* (for  $p = 2$ ) if

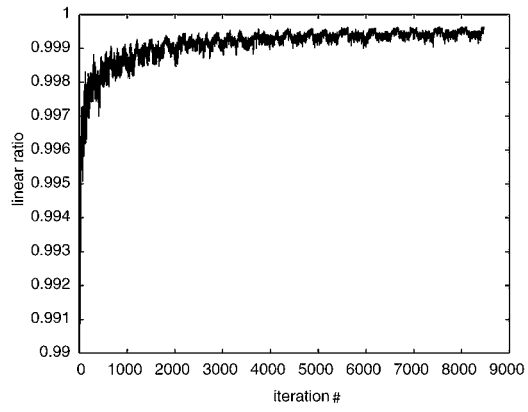
$$\frac{\|x_{k+1} - x^*\|}{\|x_k - x^*\|^p} \leq M, \quad \text{for all } k \text{ sufficiently large,}$$

where  $M$  is a positive constant not necessarily less than 1. It can be easily seen that a Q-convergent sequence of order strictly greater than 1 converges Q-super-linearly.<sup>12</sup>

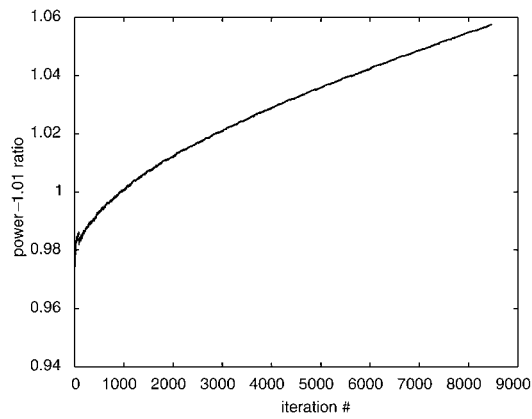
Convergence rates can be observed experimentally by recording the values of the objective function through the course of iteration and plotting the respective ratios versus the iteration number. Sample plots are shown in figures 2 and 3 for both problems on a training sample of size 10000. Limit points have been obtained from training without decomposition. Similar plots have been observed in other experiments.

It is evident from the convergence plots that the decomposition algorithm converges linearly but not super-linearly or quadratically. The plots for the linear ratio are bounded by 1, while the plots for other two ratios are unbounded.<sup>13</sup> The plots also reveal that the training problem is very ill-conditioned, as the ratio stays very close to 1.

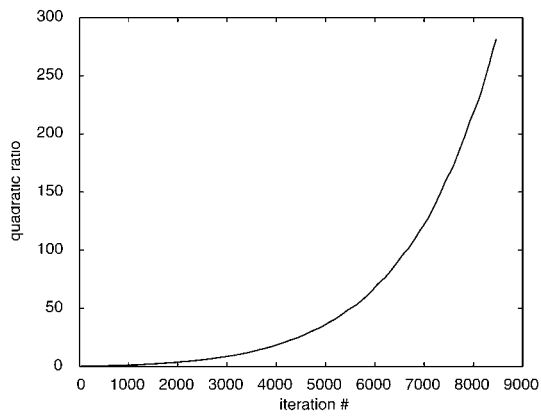
The above results are consistent with known linear convergence rates for gradient descent methods of unconstrained optimization. Noteworthy is particular semblance for the full feature set problem in which most of the variables stay away from the upper bounds and thus resemble an unconstrained case. Another important message is that convergence analysis (experimental or theoretical) is of special importance for decomposition algorithms. Unlike the scaling factors it reveals the effects of conditioning on the algorithm’s performance.



(a) linear



(b)  $p = 1.01$



(c) quadratic

Figure 2. Convergence rates for the full feature space.

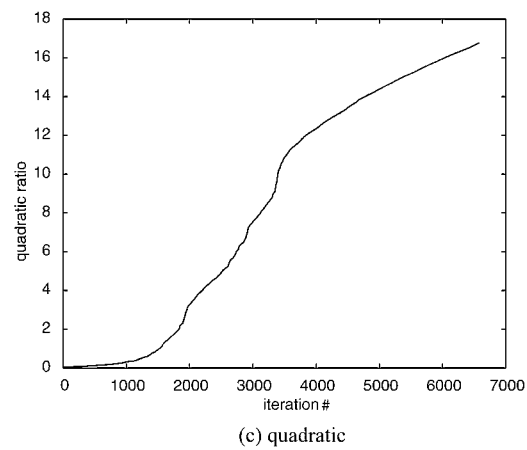
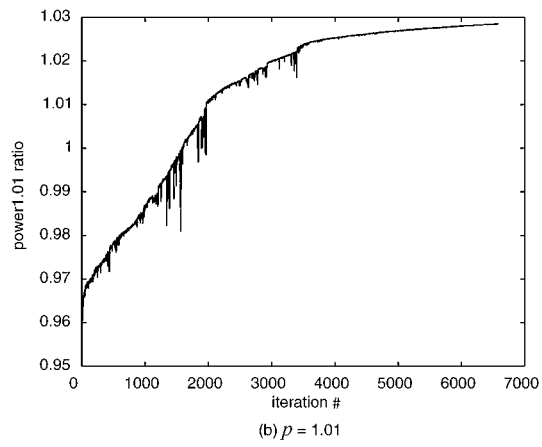
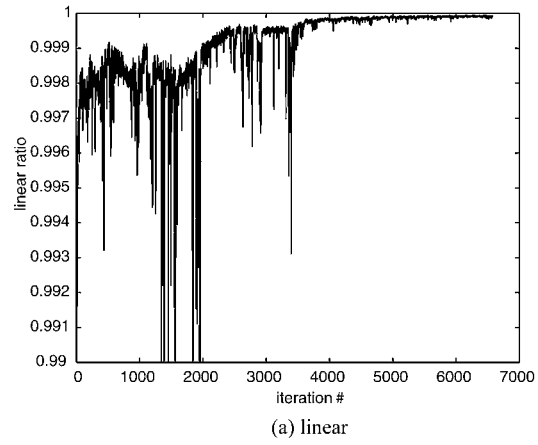


Figure 3. Convergence rates for the reduced feature space.



### 6.3. *Profiled scaling factors*

The overall processing performed during an iteration of the feasible direction decomposition algorithm can be broken up into 5 main steps:

1. Optimization: optimization proper and calculation of support vectors.
2. Update of gradients.
3. Kernel evaluation: all requests for dot products (from all modules of the system). Notice that with caching of kernel evaluations, this operation is not equally distributed across iterations: at the beginning it takes longer when kernels must be computed; towards the end all kernels end up in cache.<sup>14</sup>
4. Selection: computation of maximal violation of KKT conditions, left and right passes.
5. Evaluation: computation of objective function and threshold.

In the following section the scaling factors per iteration are established for the five factors (overall scaling factors for kernel evaluation).

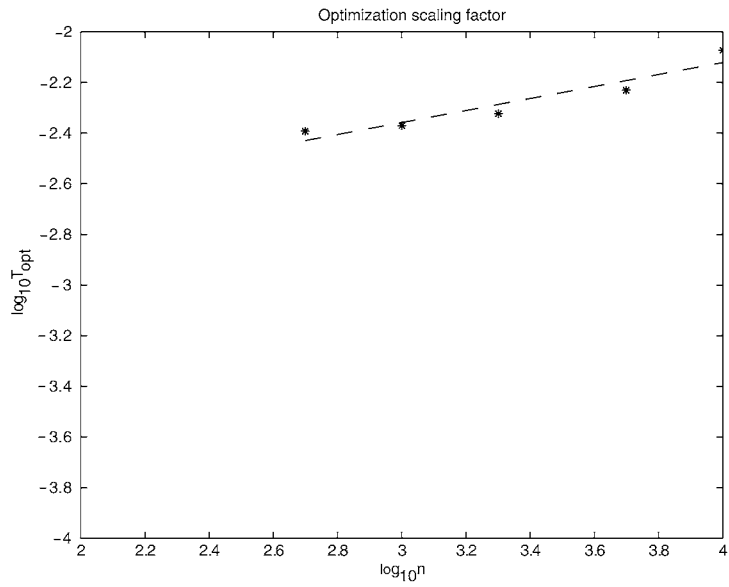
**6.3.1. *Optimization scaling factors.*** The values of obtained selection scaling factors are 0.237 for the full feature space and 0.176 for the reduced feature space, however, the quality of fits is low. The expected behavior was constant-time per iteration, because the working sets are of constant size. Perhaps, for larger data sets conditioning of sub-problems deteriorates slightly, thus increasing the number of iterations in the optimizer. The fits are displayed in figure 4.

**6.3.2. *Update scaling factors.*** The values of obtained update scaling factors are 1.060 for the full feature space and 1.064 for the reduced feature space. This coincides with the theoretical expectations of linear growth order of update operation. The fits are displayed in figure 5.

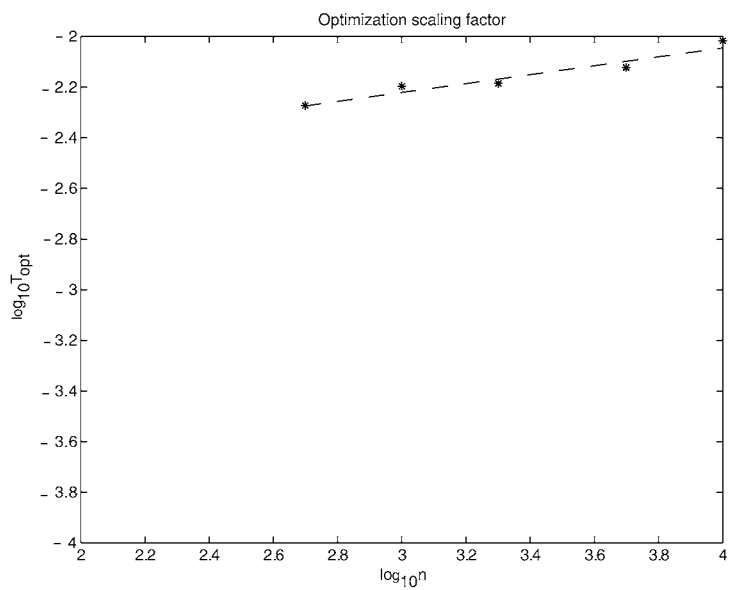
**6.3.3. *Kernel scaling factors.*** Kernel scaling factors are computed based on the timing accumulated over the entire run. The obtained values are 2.098 and 2.067 for the full and the reduced sets of features respectively. This coincides with the expected quadratic order of growth. The fits are displayed in figure 6.

**6.3.4. *Selection scaling factors.*** The values of obtained selection scaling factors are 1.149 for the full feature space and 1.124 for the reduced feature space. This is close to the theoretical expectations of linear growth order of selection operation.<sup>15</sup> The fits are displayed in figure 7.

**6.3.5. *Evaluation scaling factors.*** The values of obtained evaluation scaling factors are 1.104 for the full feature space and 1.118 for the reduced feature space. This is close to the theoretical expectations of linear growth order of evaluation operation. The fits are displayed in figure 8.



(a) full set of features



(b) reduced set of features

*Figure 4.* Optimization scaling factor fits.

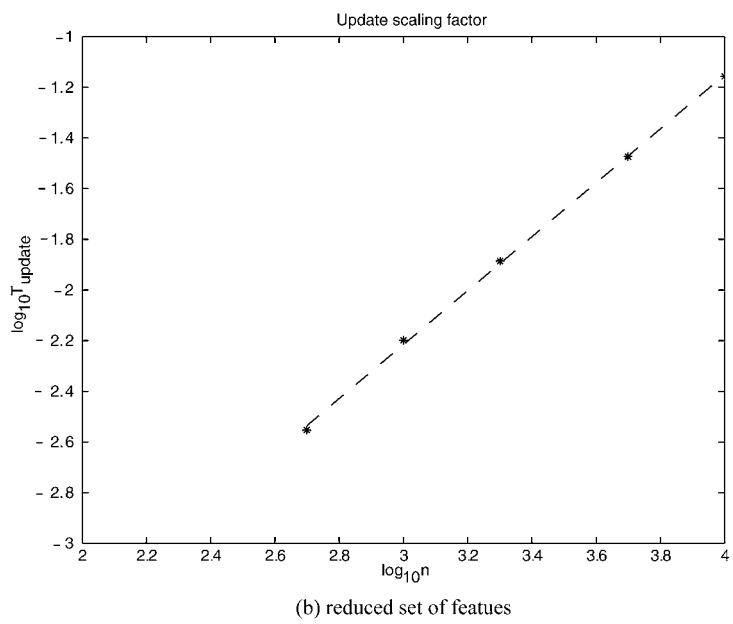
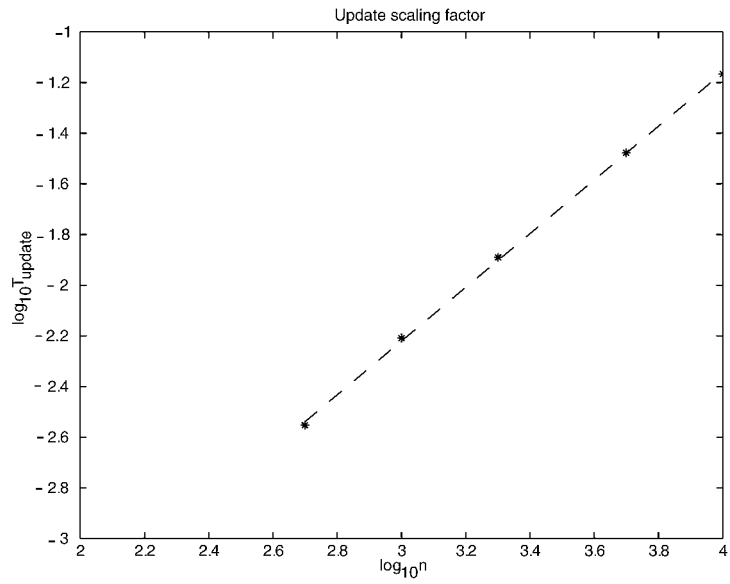
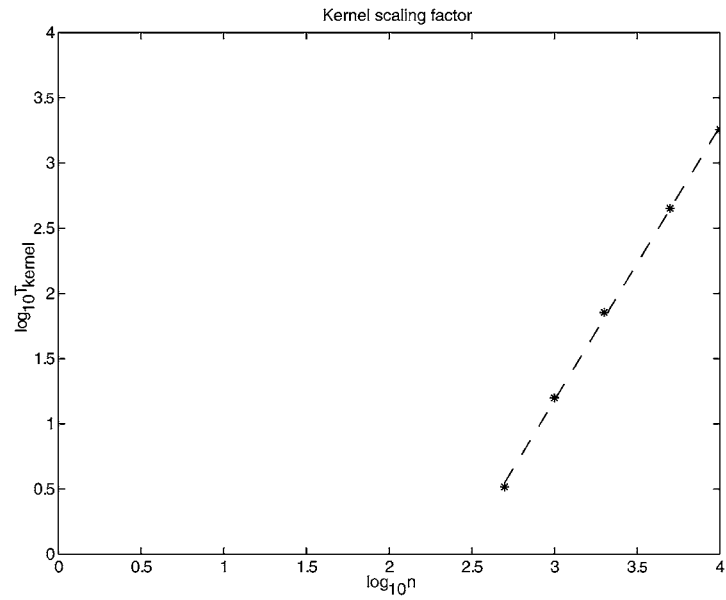
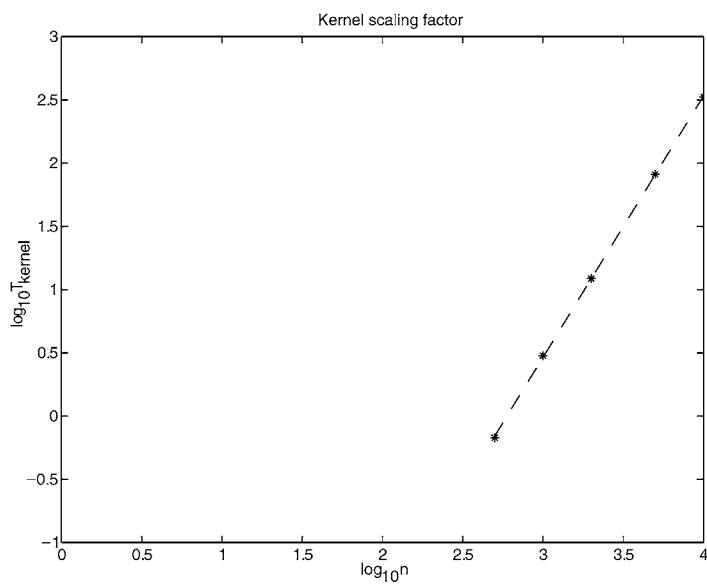


Figure 5. Update scaling factor fits.



(a) full set of features



(b) reduced set of features

*Figure 6.* Kernel scaling factor fits.

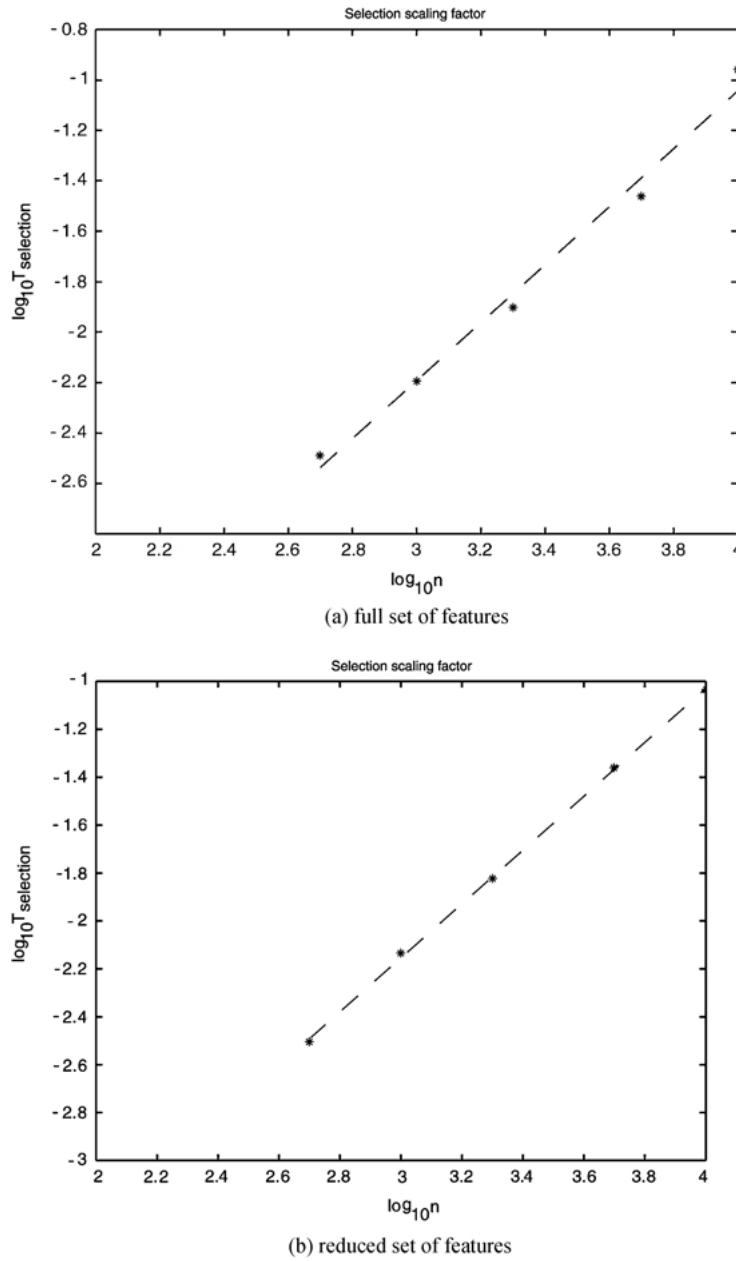
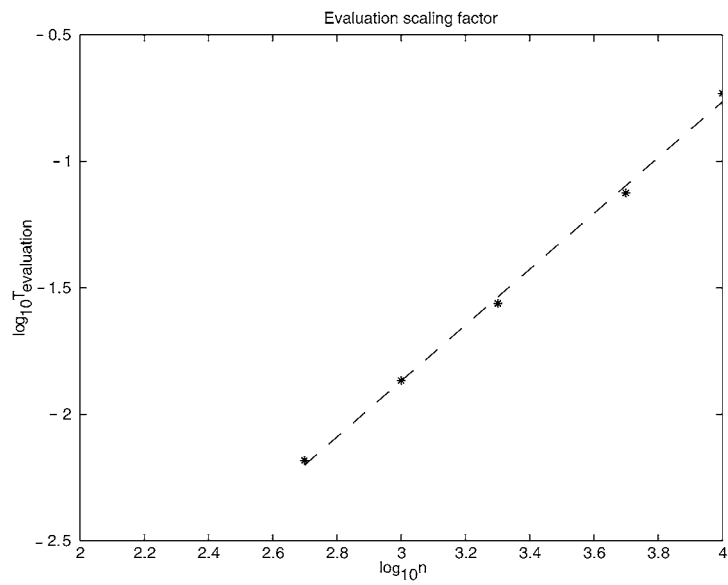
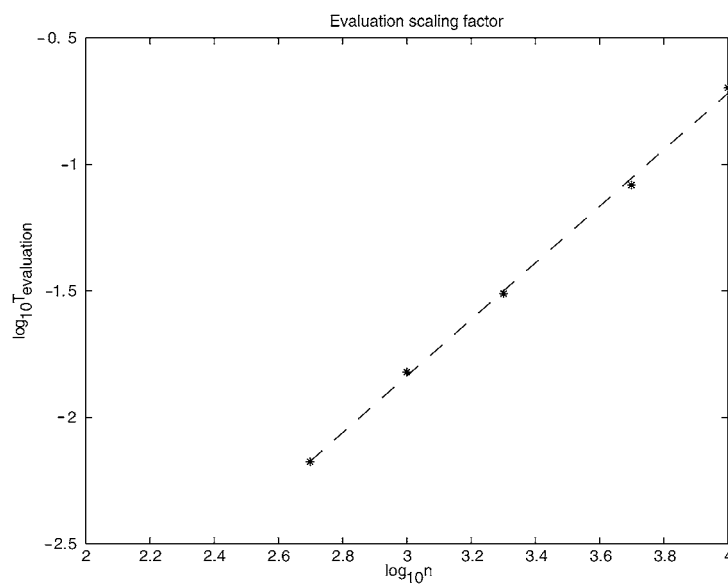


Figure 7. Selection scaling factor fits.



(a) full set of features



(b) reduced set of features

Figure 8. Evaluation scaling factor fits.

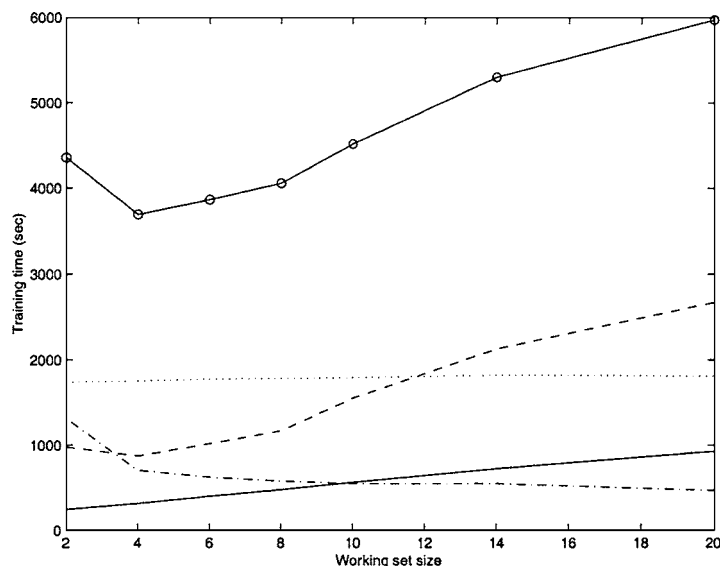


Figure 9. Training time (sec) for the KDD Cup problem for various working set sizes, full feature set, 10000 examples.

#### 6.4. Determination of working set size

The profiled timing experiments have been performed for the working sets ranging from size 2 to size 20, in the same setup as the experiments in the previous section. The results are shown in figures 9 and 10. The circled solid lines display total execution time, the dashed lines—evaluation time, dot-dashed lines—selections time, regular solid lines—update time, dotted line—kernel evaluation time. Optimization time was inferior and is not shown on the plots.

The seemingly most apparent conclusion is that the optimal working set is very small, perhaps 2 or 4. However, examination of the components of the training time reveals that this may vary for different problems and different implementations of the algorithm. While the kernel computation time is almost independent of the working set size, one can see that selection time decreases with the growth of the working set, while evaluation and update times increase. The reason for such different behavior lies in different orders of growth of these components with respect to the working set size  $q$ . The selection time is of order  $\log q$ , whereas the evaluation and the update times are of order  $q$ , per iteration. The number of iterations decreases with growing working set size, although the order of decrease is hard to establish because it depends on convergence properties of the decomposition algorithm. However, this order seems to be between logarithmic and linear—exemplified by the fact that this decrease dominates the logarithmically growing selection time but is inferior to the linearly growing update and evaluation times. One can also notice that if the three factors are more balanced<sup>16</sup> the optimal working set may be significantly larger. This fact demonstrates

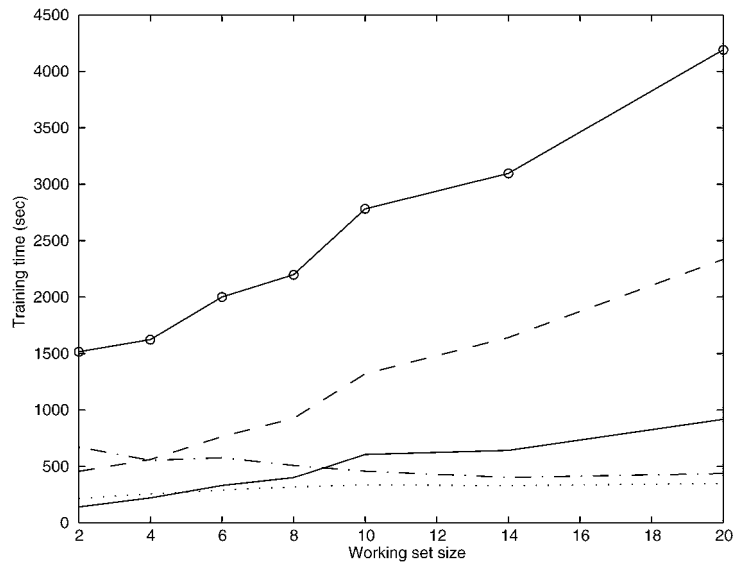


Figure 10. Training time (sec) for the KDD Cup problem for various working set sizes, reduced feature set, 10000 examples.

the possible advantage of the feasible direction decomposition over SMO: it adds another implementation-dependent leverage to improve computational efficiency.

Similar results were observed for the value of  $C = 10$ . The running times are significantly larger, because the algorithm has to “climb” in a larger bounding box of constraints, however the relationship between the contributions of the main steps of the algorithm is similar.

## 7. Discussion

To complement presentation of the feasible direction decomposition algorithms, we would like to discuss two main alternative approaches to speeding up training of SVM that have been proposed in the last years. Some of them have been mentioned in the main exposition; however, the ensuing discussion provides a more systematic review and comparison.

The most popular group of algorithms utilizes the idea of Sequential Minimal Optimization (SMO). Initiated in the work of Platt (1999), it has received a number of important extensions. Its adaptation to the regression SVM has been done in Smola and Schölkopf (1998).<sup>17</sup> A series of important results related to SMO are due to S. Keerthi and his group. In Keerthi et al. (1999) and Shevade et al. (1999) they address the threshold computation in SMO. They present examples on which SMO spends a significant amount of time at the end of training trying to establish the single optimal threshold  $b$  which is involved in the Osuna-style KKT conditions (c.f. Eq. (6) and (9)). A significant improvement in efficiency is reportedly achieved by maintaining two thresholds  $b_{\text{low}}$  and  $b_{\text{up}}$ , and using the termination condition  $b_{\text{up}} - b_{\text{low}} < \tau$  for some tolerance parameter  $\tau$ . A similar construction



involving  $\tau$  and insensitivity  $\epsilon$  is used for regression. One can easily see that these conditions are of the same flavor as the conditions of the maximal inconsistency algorithm: both approaches maintain the interval of allowable values of the Lagrangian multiplier associated with the equality constraint. Keerthi's conditions do not seem to account for a possibility of either value to be infinite, which may happen if all examples of one class (or examples lying above or below the hyperplane representing the regression function) have  $\alpha$ 's strictly at the bound. In general, the SMO-like algorithms differ from the feasible direction algorithms in that their working set selection is rather heuristic, which makes it difficult to compare these two approaches. A comprehensive numerical evaluation is highly desirable.

Another important group of algorithms builds on the idea of Successive Over-relaxation (SOR). The original algorithms for the pattern recognition and the regression SVM appeared in Mangasarian and Musicant (1999b) and Mangasarian and Musicant (1999a) respectively. This method redefines the formulation of the SVM problem to obtain the quadratic program lacking the equality constraint. As a result, the SOR method, well-known from the classical optimization literature, can be directly applied, coupled with a specially designed chunking technique. This becomes possible because, with only inequality constraints present, the values of the variables can be adjusted one at a time. A modification of this algorithm, due to Hsu and Lin (1999), deals with balanced selection of positive and negative examples and avoids alternating working sets. The numerical experiments presented in the latter work yield performance similar to that of SVM<sup>light</sup>, the implementation of Joachims' algorithm. The advantage of SOR-based methods is that they are better understood from the theoretical point of view; in particular, they are known to converge linearly (Mangasarian and Musicant, 1999b), Theorem 3.3. Their disadvantage lies in the re-formulation of the SVM training problem: it has been shown in Cristianini and Shawe-Taylor (2000) that this approach can increase VC-dimension and thus yield higher classification error. In general, it is likewise hard to draw a comparison between this group of algorithms and the feasible direction methods except by experimental evaluation.

Lastly, we would like to remark that, independently from the presented work, an identical algorithm for the regression SVM has been recently developed by utilizing the direct extension of Joachims' algorithm for the regression SVM (Collobert and Bengio, 2000). Their numerical results demonstrate significantly lower training time in comparison with SMO.

## 8. Conclusions

The unified treatment of the working set selection in SVM decomposition algorithms presented in this article provides a general view of decomposition methods rooted in the principle of feasible direction, regardless of the particular SVM formulation. The new algorithm, based on the idea of maximal inconsistency of KKT conditions, is realized in a conceptually similar fashion for both the pattern recognition and the regression SVM. Formal justification of the maximal inconsistency strategy has been presented which provides a useful insight into the mechanism of working set selection.

The experimental results presented in the article demonstrate that, similar to the pattern recognition case, significant decrease of training time can be achieved by using the decomposition algorithm. While the scaling factors of the decomposition algorithms are significantly better than those of straightforward optimization, a word of caution needs to be said with regard to the constants. It can be seen from the profiled experiments that the worst-case growth orders can only be given on a per-iteration basis, and the number of iterations, which depends on the convergence rate and the required precision, adds another dimension to running time analysis. An important experimental result establishes the linear convergence rate of the feasible direction decomposition algorithms. It was also observed that problem conditioning can make significant impact on the constants involved in linear convergence.

A number of open questions remain regarding the SVM decomposition algorithms. Can the linear convergence rate be established theoretically? Can a super-linear or a quadratic convergence rate be achieved by a different algorithm? Finally, extremely far-reaching results can be borne by investigation of conditioning of the training problem: since the latter is a by-product of a number of factors, such as the choice of kernel, the box constraint, the risk functional, etc, conditioning of the optimization problem might be useful to guide the selection of SVM parameters.

## Notes

1. Another possible bottleneck is computation of the kernel matrix. If real dimension of data points is large, this time may become comparable with training time. However it will be assumed in the rest of this article that kernel computation is not the main factor contributing to the complexity of training.
2. In the notation of Osuna, Freund, & Girosi (1997) the quantity below is denoted by  $g_i$ ; however, it makes more sense to reserve notation  $\mathbf{g}$  for the gradient of the objective function used later in this article.
3. Several alternative statements of the direction finding problem can be found in Zoutendijk (1960).
4. Recall that for the pattern recognition SVM the equality constraint is  $\mathbf{y}^T \mathbf{d} = 0$ .
5. The motivation for the names “forward” and “backward” will be clear shortly.
6. In practice, sorting which takes  $O(n \log n)$  operations can be replaced with heap-based algorithms yielding the complexity of  $O(n \log q)$  or  $O(q \log n)$  depending on how the heap is built.
7. They will be in the regression case.
8. Since both algorithms use identical feasibility check for every sample, it is obvious that infeasible samples will never be included in a working set in both algorithms.
9. In order for the results to be conceptually compatible with Joachims’ I used the old point-wise termination conditions.
10. Without decomposition, MINOS was used to solve the full-size problem.
11. The scaling factors also vary in Joachims (1999).
12. The prefix “Q” will be omitted in the rest of the presentation in this section.
13. It is difficult to see from the plot for the reduced feature set, but there is a tiny margin of  $\sim 0.0001$  which separates the ratios from 1.
14. In the experiments above enough memory was allocated for the kernel cache to hold the values of kernels for all support vectors.
15. The implementation uses heap-based method whose theoretical running time order is  $O(n \log q)$ . The logarithmic factor does not feature in the scaling factor because  $q$  is assumed constant.
16. For example, the dominant evaluation time can be readily eliminated because computation of objective function at every iteration is not crucial for decomposition and computation of the threshold can be avoided in the maximal inconsistency algorithm; in the current implementation it was purposely done in the same way as in Joachims’ algorithm.
17. Platt’s work was available as a technical report in 1998, hence the seemingly bizarre reference.

## References

- Chang, C.-C., Hsu, C.-W., & Lin, C.-J. (1999). The analysis of decomposition methods for support vector machines. In *Proceedings of the International Joint Conference on Artificial Intelligence. SVM Workshop*.
- Collobert, R. & Bengio, S. (2000). Support vector machines for large-scale regression problems. Technical Report IDIAP-RR 00-17, IDIAP.
- Cristianini, N. & Shawe-Taylor, J. (2000). *An introduction to support vector machines*. Cambridge, UK: Cambridge University Press.
- Hsu, C.-W. & Lin, C.-J. (1999). A simple decomposition method for support vector machines. Technical report, National Taiwan University.
- Joachims, T. (1999). Making large-scale support vector machine learning practical. In Schölkopf et al. (Eds.). *Advances in Kernel Methods-Support Vector Learning* (pp. 169–184). Cambridge, MA: MIT Press.
- Kaufmann, L. (1999). Solving the quadratic problem arising in support vector classification. In Schölkopf et al. (Eds.). *Advances in Kernel Methods-Support Vector Learning* (pp. 147–168). Cambridge, MA: MIT Press.
- Keerthi, S. S., Shevade, S. K., Bhattacharyya, C., & Murthy, K. R. K. (1999). Improvements to Platt's SMO algorithm for SVM classifier design. Technical Report CD-99-14, National Institute of Singapore.
- Laskov, P. (2000). An improved decomposition algorithm for regression support vector machines. In S. Solla, T. Leen, & K.-R. Müller (Eds.). *Advances in neural information processing systems 12* (pp. 484–490). Cambridge, MA: MIT Press.
- Mangasarian, O. L. (1969). *Nonlinear programming*. New York: McGraw-Hill.
- Mangasarian, O. L. & Musicant, D. R. (1999a). Massive support vector regression. Technical Report 99-02, University of Wisconsin, Data Mining Institute.
- Mangasarian, O. L. & Musicant, D. R. (1999b). Successive overrelaxation for support vector machines. *Transactions on Neural Networks*, 10:5, 100–106.
- Nocedal, J. & Wright, S. J. (1999). *Numerical optimization*. New York, NY: Springer-Verlag.
- Osuna, E. (1998). Support vector machines: Training and applications. Ph.D. Thesis, Massachusetts Institute of Technology.
- Osuna, E., Freund, R., & Girosi, F. (1997). Improved training algorithm for support vector machines. *Neural Networks and Signal Processing*.
- Platt, J. C. (1999). Fast training of support vector machines using sequential minimal optimization. In Schölkopf et al. (Eds.) 1999. *Advances in kernel method-support vector training* (pp. 185–208). Cambridge, MA: MIT Press.
- Schölkopf B. (1997). Support vector learning. Ph.D. Thesis, Universität Tübingen.
- Schölkopf B., Burges, C., & Smola, A. (Eds.) (1999). *Advances in kernel methods—support vector learning*. Cambridge, MA: MIT Press.
- Shevade, S. K., Keerthi, S. S., Bhattacharyya, C., & Murthy, K. R. K. (1999). Improvements to SMO algorithm for SVM regression. Technical Report CD-99-16, National Institute of Singapore.
- Smola, A. (1998). Learning with Kernels. Ph.D. Thesis, Technical University of Berlin.
- Smola, A. & Schölkopf B. (1998). A tutorial on support vector regression. Technical Report NC2-TR-1998-030, NeuroCOLT2.
- Vapnik, V. N. (1982). *Estimation of dependences based on empirical data*. New York, NY: Springer-Verlag.
- Vapnik, V. N. (1995). *The nature of statistical learning theory*. New York, NY: Springer-Verlag.
- Vapnik, V. N. (1999). *Statistical learning theory*. New York, NY: J. Wiley and Sons.
- Zoutendijk, G. (1960). *Methods of feasible directions*. Amsterdam: Elsevier.

Received Mar 14, 2000

Revised Mar 2, 2001

Accepted Mar 2, 2001

Final manuscript July 20, 2001