

## Improvements to Platt's SMO Algorithm for SVM Classifier Design

**S. S. Keerthi**

*Department of Mechanical and Production Engineering, National University of Singapore, Singapore-119260*

**S. K. Shevade**

**C. Bhattacharyya**

**K. R. K. Murthy**

*Department of Computer Science and Automation, Indian Institute of Science, Bangalore-560012, India*

**This article points out an important source of inefficiency in Platt's sequential minimal optimization (SMO) algorithm that is caused by the use of a single threshold value. Using clues from the KKT conditions for the dual problem, two threshold parameters are employed to derive modifications of SMO. These modified algorithms perform significantly faster than the original SMO on all benchmark data sets tried.**

### 1 Introduction ---

In the past few years, there has been a lot of excitement and interest in support vector machines (Vapnik, 1995; Burges, 1998) because they have yielded excellent generalization performance on a wide range of problems. Recently, fast iterative algorithms that are also easy to implement have been suggested (Platt, 1998; Joachims, 1998; Mangasarian & Musicant, 1998; Friess, 1998; Keerthi, Shevade, Bhattacharyya, & Murthy, 1999a). Platt's sequential minimization algorithm (SMO) (Platt, 1998, 1999) is an important example. A remarkable feature of SMO is that it is also extremely easy to implement. Platt's comparative testing against other algorithms has shown that SMO is often much faster and has better scaling properties.

In this article we enhance the value of SMO even further. In particular, we point out an important source of inefficiency caused by the way SMO maintains and updates a single threshold value. Getting clues from optimality criteria associated with the KKT conditions for the dual, we suggest the use of two threshold parameters and devise two modified versions of SMO that are more efficient than the original SMO. Computational comparison on benchmark data sets shows that the modifications perform significantly faster than the original SMO in most situations. The ideas mentioned in this article can also be applied to the SMO regression algorithm (Smola, 1998).

We report the results of that extension in Shevade, Keerthi, Bhattacharyya, and Murthy (1999).

In section 2 we briefly discuss the SVM problem formulation, the dual problem, and the associated KKT optimality conditions. We also point out how these conditions lead to proper criteria for terminating algorithms for designing SVM classifiers. Section 3 gives a short summary of Platt's SMO algorithm. In section 4 we point out the inefficiency associated with the way SMO uses a single threshold value, and we describe the modified algorithms in section 5. Computational comparison is done in section 6.

## 2 SVM Problem and Optimality Conditions

The basic problem addressed in this article is the two-category classification problem. Burges (1998) gives a good overview of the solution of this problem using SVMs. Throughout this article we will use  $x$  to denote the input vector of the SVM and  $z$  to denote the feature space vector, which is related to  $x$  by a transformation,  $z = \phi(x)$ . As in all other SVM designs, we do not assume  $\phi$  to be known; all computations will be done using only the kernel function,  $k(x, \hat{x}) = \phi(x) \cdot \phi(\hat{x})$ , where  $\cdot$  denotes inner product in the  $z$  space. Let  $\{(x_i, y_i)\}$  denote the training set, where  $x_i$  is the  $i$ th input pattern and  $y_i$  is the corresponding target value;  $y_i = 1$  means  $x_i$  is in class 1, and  $y_i = -1$  means  $x_i$  is in class 2. Let  $z_i = \phi(x_i)$ . The optimization problem solved by the SVM is:

$$\min \frac{1}{2} \|w\|^2 + C \sum_i \xi_i \quad \text{s.t.} \quad y_i(w \cdot z_i - b) \geq 1 - \xi_i \quad \forall i; \quad \xi_i \geq 0 \quad \forall i. \quad (2.1)$$

This problem is referred to as the primal problem.

Let us define  $w(\alpha) = \sum_i \alpha_i y_i z_i$ . We will refer to the  $\alpha_i$ 's as Lagrange multipliers. The  $\alpha_i$ 's are obtained by solving the following dual problem:

$$\max W(\alpha) = \sum_i \alpha_i - \frac{1}{2} w(\alpha) \cdot w(\alpha) \quad \text{s.t.} \quad 0 \leq \alpha_i \leq C \quad \forall i; \quad \sum_i \alpha_i y_i = 0. \quad (2.2)$$

Once the  $\alpha_i$ 's are obtained, the other primal variables,  $w$ ,  $b$ , and  $\xi$ , can be easily determined by using the KKT conditions for the primal problem. It is possible that the solution is nonunique; for instance, when all  $\alpha_i$ s take the boundary values of 0 and  $C$ , it is possible that  $b$  is not unique.

The numerical approach in SVM design is to solve the dual (instead of the primal) because it is a finite-dimensional optimization problem. (Note that  $w(\alpha) \cdot w(\alpha) = \sum_i \sum_j y_i y_j \alpha_i \alpha_j k(x_i, x_j)$ .) To derive proper stopping conditions for algorithms that solve the dual, it is important to write down the optimality conditions for the dual. The Lagrangian for the dual is:

$$\tilde{L} = \frac{1}{2} w(\alpha) \cdot w(\alpha) - \sum_i \alpha_i - \sum_i \delta_i \alpha_i + \sum_i \mu_i (\alpha_i - C) - \beta \sum_i \alpha_i y_i.$$

Define

$$F_i = w(\alpha) \cdot z_i - y_i = \sum_j \alpha_j y_j k(x_i, x_j) - y_i.$$

The KKT conditions for the dual problem are:

$$\frac{\partial \bar{L}}{\partial \alpha_i} = (F_i - \beta)y_i - \delta_i + \mu_i = 0, \quad \delta_i \geq 0, \quad \delta_i \alpha_i = 0, \quad \mu_i \geq 0, \quad \mu_i(\alpha_i - C) = 0 \quad \forall i.$$

These conditions<sup>1</sup> can be simplified by considering three cases for each  $i$ :

Case 1.  $\alpha_i = 0$ :

$$\delta_i \geq 0, \quad \mu_i = 0 \Rightarrow (F_i - \beta)y_i \geq 0. \quad (2.3a)$$

Case 2.  $0 < \alpha_i < C$ :

$$\delta_i = 0, \quad \mu_i = 0 \Rightarrow (F_i - \beta)y_i = 0. \quad (2.3b)$$

Case 3.  $\alpha_i = C$ :

$$\delta_i = 0, \quad \mu_i \geq 0 \Rightarrow (F_i - \beta)y_i \leq 0. \quad (2.3c)$$

Define the following index sets at a given  $\alpha$ :  $I_0 = \{i: 0 < \alpha_i < C\}$ ;  $I_1 = \{i: y_i = 1, \alpha_i = 0\}$ ;  $I_2 = \{i: y_i = -1, \alpha_i = C\}$ ;  $I_3 = \{i: y_i = 1, \alpha_i = C\}$ ; and,  $I_4 = \{i: y_i = -1, \alpha_i = 0\}$ . Note that these index sets depend on  $\alpha$ . The conditions in equations 2.3a through 2.3c can be rewritten as

$$\beta \leq F_i \quad \forall i \in I_0 \cup I_1 \cup I_2; \quad \beta \geq F_i \quad \forall i \in I_0 \cup I_3 \cup I_4. \quad (2.4)$$

It is easily seen that optimality conditions will hold iff there exists a  $\beta$  satisfying equation 2.4. The following result is an immediate consequence.

**Lemma 1.** *Define:*

$$b_{\text{up}} = \min\{F_i: i \in I_0 \cup I_1 \cup I_2\} \quad \text{and} \quad b_{\text{low}} = \max\{F_i: i \in I_0 \cup I_3 \cup I_4\}. \quad (2.5)$$

*Then optimality conditions will hold at some  $\alpha$  iff*

$$b_{\text{low}} \leq b_{\text{up}}. \quad (2.6)$$

---

<sup>1</sup> The KKT conditions are both necessary and sufficient for optimality. Hereafter we will simply refer to them as optimality conditions.

It is easy to see the close relationship between the threshold parameter  $b$  in the primal problem and the multiplier,  $\beta$ . In particular, at optimality,  $\beta$  and  $b$  are identical. Therefore, in the rest of the article,  $\beta$  and  $b$  will denote the same quantity.

We will say that an index pair  $(i, j)$  defines a *violation* at  $\alpha$  if *one* of the following sets of conditions holds:

$$i \in I_0 \cup I_3 \cup I_4, \quad j \in I_0 \cup I_1 \cup I_2 \quad \text{and} \quad F_i > F_j \quad (2.7a)$$

$$i \in I_0 \cup I_1 \cup I_2, \quad j \in I_0 \cup I_3 \cup I_4 \quad \text{and} \quad F_i < F_j. \quad (2.7b)$$

Note that optimality conditions will hold at  $\alpha$  iff there does not exist any index pair  $(i, j)$  that defines a violation.

Since, in numerical solution, it is usually not possible to achieve optimality exactly, there is a need to define approximate optimality conditions. Condition 2.6 can be replaced by

$$b_{\text{low}} \leq b_{\text{up}} + 2\tau, \quad (2.8)$$

where  $\tau$  is a positive tolerance parameter. (In the pseudocode given in Platt, 1998, this parameter is referred to as  $\text{tol}$ .) Correspondingly, the definition of violation can be altered by replacing equations 2.7a and 2.7b by:

$$i \in I_0 \cup I_3 \cup I_4, \quad j \in I_0 \cup I_1 \cup I_2 \quad \text{and} \quad F_i > F_j + 2\tau \quad (2.9a)$$

$$i \in I_0 \cup I_1 \cup I_2, \quad j \in I_0 \cup I_3 \cup I_4 \quad \text{and} \quad F_i < F_j - 2\tau. \quad (2.9b)$$

Hereafter, when optimality is mentioned, it will mean approximate optimality.

Since  $\beta$  can be placed halfway between  $b_{\text{low}}$  and  $b_{\text{up}}$ , approximate optimality conditions will hold iff there exists a  $\beta$  such that equations 2.3a through 2.3c are satisfied with a  $\tau$ -margin, that is:

$$(F_i - \beta)y_i \geq -\tau \quad \text{if} \quad \alpha_i = 0 \quad (2.10a)$$

$$|(F_i - \beta)| \leq \tau \quad \text{if} \quad 0 < \alpha_i < C \quad (2.10b)$$

$$(F_i - \beta)y_i \leq \tau \quad \text{if} \quad \alpha_i = C. \quad (2.10c)$$

These three equations, 2.10a through 2.10c, are the approximate optimality conditions employed by Platt (1998), Joachims (1998), and others. Note that if equations 2.10a through 2.10c are violated for one value of  $\beta$ , that does not mean that optimality conditions are violated.

It is also possible to give an alternate approximate optimality condition based on the closeness of  $W(\alpha)$  to the optimal value, estimated using duality gap ideas. This is fine, but care is needed in choosing the tolerance used (see Keerthi et al., 1999a, for a related discussion). This criterion has the disadvantage that it is a single global condition involving all  $i$ . On the

other hand, equation 2.9 consists of an individual condition for each pair of indices, and hence it is much better suited to the methods discussed here. In particular, when  $(i, j)$  satisfies equation 2.9, then a strict improvement in the dual objective function can be achieved by optimizing only  $\alpha_i$  and  $\alpha_j$ . (This is true even if  $\tau = 0$ .)

### 3 Platt's SMO Algorithm

---

We now give a brief description of the SMO algorithm. The basic step consists of choosing a pair of indices,  $(i_1, i_2)$  and optimizing the dual objective function by adjusting  $\alpha_{i_1}$  and  $\alpha_{i_2}$  only. Because the working set is only of size 2 and the equality constraint in equation 2.2 can be used to eliminate one of the two Lagrange multipliers, the optimization problem at each step is a quadratic minimization in just one variable. It is straightforward to write down an analytic solution for it. (Complete details are given in Platt, 1998.) The procedure `takeStep` (which is a part of the pseudocode given there) gives a clear description of the implementation. There is no need to recall all details here. We make only one important comment on the role of the threshold parameter,  $\beta$ . As in Platt (1998) define the output error on the  $i$ th pattern as

$$E_i = F_i - \beta.$$

Consistent with the pseudocode of Platt (1998), let us call the indices of the two multipliers chosen for optimization in one step as  $i_2$  and  $i_1$ . A look at the details in Platt (1998) shows that to take a step by varying  $\alpha_{i_1}$  and  $\alpha_{i_2}$ , we only need to know  $E_{i_1} - E_{i_2} = F_{i_1} - F_{i_2}$ . Therefore, knowledge of the value of  $\beta$  is not needed to take a step.

The method followed to choose  $i_1$  and  $i_2$  at each step is crucial for efficient solution of the problem. Based on a number of experiments, Platt came up with a good set of heuristics. He employs a two-loop approach: the outer loop chooses  $i_2$ , and for a chosen  $i_2$ , the inner loop chooses  $i_1$ . The outer loop iterates over all patterns, violating the optimality conditions—first only over those with Lagrange multipliers on neither the upper nor lower boundary (in Platt's pseudocode, this looping is indicated by `examineAll=0`) and, once all of them are satisfied, over all patterns violating the optimality conditions (`examineAll=1`) to ensure that the problem has indeed been solved. For efficient implementation Platt maintains and updates a cache for  $E_i$  values for indices  $i$  corresponding to nonboundary multipliers. The remaining  $E_i$  are computed as and when needed.

Let us now see how the SMO algorithm chooses  $i_1$ . The aim is to make a large increase in the objective function. Since it is expensive to try out all possible choices of  $i_1$  and choose the one that gives the best increase in objective function, the index  $i_1$  is chosen to maximize  $|E_{i_2} - E_{i_1}|$ . (If we define  $\rho(t) = W(\alpha(t))$  where  $\alpha_{i_1}(t) = \alpha_{i_1}^{\text{old}} + y_{i_1}t$ ,  $\alpha_{i_2}(t) = \alpha_{i_2}^{\text{old}} - y_{i_2}t$ , and  $\alpha_i = \alpha_i^{\text{old}} \forall i \notin \{i_1, i_2\}$ , then  $|\rho'(0)| = |E_{i_1} - E_{i_2}|$ .) Since  $E_{i_2}$  is available in cache for

nonboundary multiplier indices, only such indices are initially used in the above choice of  $i_1$ . If such a choice of  $i_1$  does not yield sufficient progress, then the following steps are taken. Starting from a randomly chosen index, all indices corresponding to nonbound multipliers are tried, one by one, as choices for  $i_1$ . If sufficient progress still is not possible, all indices are tried, one by one, as choices for  $i_1$ , again starting from a randomly chosen index. Thus, the choice of random seed affects the running time of SMO.

Although a value of  $\beta$  is not needed to take a step, it is needed if equations 2.10a through 2.10c are employed for checking optimality. In the SMO algorithm,  $\beta$  is updated after each step. A value of  $\beta$  is chosen so as to satisfy equation 2.3 for  $i \in \{i_1, i_2\}$ . If, after a step involving  $(i_1, i_2)$ , one of  $\alpha_{i_1}$ ,  $\alpha_{i_2}$  (or both) takes a nonboundary value, then equation 2.3b is exploited to update the value of  $\beta$ . In the rare case that this does not happen, there exists a whole interval, say,  $[\beta_{\text{low}}, \beta_{\text{up}}]$ , of admissible thresholds. In this situation SMO simply chooses:  $\beta = (\beta_{\text{low}} + \beta_{\text{up}})/2$ .

#### 4 Inefficiency of the SMO Algorithm

SMO is a carefully organized algorithm with excellent computational efficiency. However, because of its way of computing and using a single threshold value, it can carry some inefficiency unnecessarily. At any instant, the SMO algorithm fixes  $\beta$  based on the current two indices that are being optimized. However, while checking whether the remaining examples violate optimality, it is quite possible that a different, shifted choice of  $\beta$  may do a better job. So in the SMO algorithm, it is quite possible that although  $\alpha$  has reached a value where optimality is satisfied (that is, equation 2.8), SMO has not detected this because it has not identified the correct choice of  $\beta$ . It is also quite possible that a particular index may appear to violate the optimality conditions because equation 2.10 is employed using an “incorrect” value of  $\beta$ , although this index may not be able to pair with another to define a violation. In such a situation, the SMO algorithm does an expensive and wasteful search looking for a second index so as to take a step. We believe that this is an important source of inefficiency in the SMO algorithm.

There is one simple alternate way of choosing  $\beta$  that involves all indices. By duality theory, the objective function value in equation 2.1 of a primal feasible solution is greater than or equal to the objective function value in equation 2.2 of a dual feasible solution. The difference between these two values is referred to as the duality gap. The duality gap is zero only at optimality. Suppose  $\alpha$  is given and  $w = w(\alpha)$ . The  $\xi_i$  can be chosen optimally (as a function of  $\beta$ ). The result is that the duality gap is expressed as a function of  $\beta$  only. One possible way of improving the SMO algorithm is always to choose  $\beta$  so as to minimize the duality gap.<sup>2</sup> This corresponds to

<sup>2</sup> One of the reviewers suggested this idea.

the subproblem,

$$\min \sum_i \max\{0, y_i(\beta - F_i)\}.$$

This is an easy problem to solve. Let  $p$  denote the number of indices in class 2, that is, the number of  $i$  having  $y_i = -1$ . In an increasing order arrangement of  $\{F_i\}$ , let  $f_p$  and  $f_{p+1}$  be the  $p$ th and  $(p + 1)$ th values. Then any  $\beta$  in the interval,  $[f_p, f_{p+1}]$  is a minimizer. The determination of  $f_p$  and  $f_{p+1}$  can be done efficiently using a median-finding technique. Since all  $F_i$  are typically not available at a given stage of the algorithm, it is appropriate to apply the above idea to that subset of indices for which  $F_i$  are available. (This set contains  $I_0$ .) We implemented this idea and tested it on some benchmark problems. It gave a mixed performance, performing well in some situations, poorly in some, and failing in a few cases. More detailed study is needed to understand the cause of this. (See section 6 for details of the performance on three examples.)

## 5 Modifications of the SMO Algorithm

In this section we suggest two modified versions of the SMO algorithm, each of which overcomes the problems mentioned in the last section. As we will see in the computational evaluation of section 6, these modifications are almost always better than the original SMO algorithm, and in most situations, they give quite good improvement in efficiency when tested on several benchmark problems.

In short, the modifications avoid the use of a single threshold value  $\beta$  and the use of equations 2.10a through 2.10c for checking optimality. Instead, two threshold parameters,  $b_{\text{up}}$  and  $b_{\text{low}}$ , are maintained, and equation 2.8 or 2.9 is employed for checking optimality. Assuming that the reader is familiar with Platt (1998) and the pseudocodes given there, we give only a set of pointers that describe the changes that are made to Platt's SMO algorithm. Pseudocodes that fully describe these changes can be found in Keerthi et al. (1999b).

1. Suppose, at any instant,  $F_i$  is available for all  $i$ . Let  $i_{\text{low}}$  and  $i_{\text{up}}$  be indices such that

$$F_{i_{\text{low}}} = b_{\text{low}} = \max\{F_i: i \in I_0 \cup I_3 \cup I_4\} \quad (5.1a)$$

$$F_{i_{\text{up}}} = b_{\text{up}} = \min\{F_i: i \in I_0 \cup I_1 \cup I_2\}. \quad (5.1b)$$

Then checking a particular  $i$  for optimality is easy. For example, suppose  $i \in I_1 \cup I_2$ . We have to check only if  $F_i < F_{i_{\text{low}}} - 2\tau$ . If this condition holds, then there is a violation, and in that case SMO's `takeStep` procedure can be applied to the index pair,  $(i, i_{\text{low}})$ . Similar steps can be given for indices in the other sets. Thus, in our approach, the checking of optimality of  $i_2$  and the choice of the second index  $i_1$  go hand in hand, unlike the original

SMO algorithm. As we will see below, we compute and use  $(i_{low}, b_{low})$  and  $(i_{up}, b_{up})$  via an efficient updating process.

2. To be efficient, we would, as in the SMO algorithm, spend much of the effort altering  $\alpha_i$ ,  $i \in I_0$ ; cache for  $F_i$ ,  $i \in I_0$  are maintained and updated to do this efficiently. And when optimality holds for all  $i \in I_0$ , only then examine all indices for optimality.

3. Some extra steps are added to the `takeStep` procedure. After a successful step using a pair of indices,  $(i_2, i_1)$ , let  $\tilde{I} = I_0 \cup \{i_1, i_2\}$ . We compute, partially,  $(i_{low}, b_{low})$  and  $(i_{up}, b_{up})$  using  $\tilde{I}$  only (that is, using only  $i \in \tilde{I}$  in equations 5.1a and 5.1b). Note that these extra steps are inexpensive because cache for  $\{F_i, i \in I_0\}$  is available and updates of  $F_{i_1}, F_{i_2}$  are easily done. A careful look shows that since  $i_2$  and  $i_1$  have been just involved in a successful step, each of the two sets,  $\tilde{I} \cap (I_0 \cup I_1 \cup I_2)$  and  $\tilde{I} \cap (I_0 \cup I_3 \cup I_4)$ , is nonempty; hence the partially computed  $(i_{low}, b_{low})$  and  $(i_{up}, b_{up})$  will not be null elements. Since  $i_{low}$  and  $i_{up}$  could take values from  $\{i_2, i_1\}$  and they are used as choices for  $i_1$  in the subsequent step (see item 1 above), we keep the values of  $F_{i_1}$  and  $F_{i_2}$  also in cache.

4. When working only with  $\alpha_i$ ,  $i \in I_0$ , that is, a loop with `examineAll=0`, one should note that if equation 2.8 holds at some point, then it implies that optimality holds as far as  $I_0$  is concerned. (This is because as mentioned in item 3 above, the choice of  $b_{low}$  and  $b_{up}$  is influenced by all indices in  $I_0$ .) This gives an easy way of exiting this loop.

5. There are two ways of implementing the loop involving indices in  $I_0$  only (`examineAll=0`):

**Method 1.** This is in line with what is done in SMO. Loop through all  $i_2 \in I_0$ . For each  $i_2$ , check optimality and, if violated, choose  $i_1$  appropriately. For example, if  $F_{i_2} < F_{i_{low}} - 2\tau$ , then there is a violation, and in that case choose  $i_1 = i_{low}$ .

**Method 2.** Always work with the worst violating pair; choose  $i_2 = i_{low}$  and  $i_1 = i_{up}$ .

Depending on which one of these methods is used, we refer to the resulting overall modification of SMO as SMO-Modification 1 and SMO-Modification 2. SMO and SMO-Modification 1 are identical except in the way optimality is tested. SMO-Modification 2 can be thought of as a further improvement of SMO-Modification 1, where the cache is effectively used to choose the violating pair when `examineAll=0`.

6. When optimality on  $I_0$  holds, we come back to check optimality on all indices (`examineAll=1`). Here we loop through all indices, one by one. Since  $(b_{low}, i_{low})$  and  $(b_{up}, i_{up})$  have been partially computed using  $I_0$  only, we update these quantities as each  $i$  is examined. For a given  $i$ ,  $F_i$  is computed first, and optimality is checked using the current  $(b_{low}, i_{low})$  and  $(b_{up}, i_{up})$ ; if there is no violation,  $F_i$  is used to update these quantities. For example, if  $i \in I_1 \cup I_2$  and  $F_i < b_{low} - 2\tau$ , then there is a violation, in which case we



Table 1: Data Set Properties.

Data Set	$\sigma^2$	$n$	$m$
Wisconsin Breast Cancer	4.0	9	683
Adult-7	10.0	123	16,100
Web-7	10.0	300	24,692

take a step using  $(i, i_{low})$ . On the other hand, if there is no violation, then  $(i_{up}, b_{up})$  is modified using  $F_i$ , that is, if  $F_i < b_{up}$  then we do:  $i_{up} := i$  and  $b_{up} := F_i$ .

7. Suppose we do as described in the previous item. What happens if there is no violation for any  $i$  in a loop having `examineAll=1`? Can we conclude that optimality holds for all  $i$ ? The answer is yes. This is easy to see from the following argument. Suppose, by contradiction, there does exist one  $(i, j)$  pair such that they define a violation, that is, they satisfy equation 2.9. Let us say  $i < j$ . Then  $j$  would not have satisfied the optimality check in the described implementation because  $F_i$  would have, earlier than  $j$  is seen, affected either the calculation of  $b_{low}$  or  $b_{up}$  settings, or both. In other words, even if  $i$  is mistakenly taken as having satisfied optimality earlier in the loop,  $j$  will be detected as violating optimality when it is analyzed. Only when equation 2.8 holds is it possible for all indices to satisfy the optimality checks. Furthermore, when equation 2.8 holds and the loop over all indices has been completed, the true values of  $b_{up}$  and  $b_{low}$ , as defined in equation 2.5, would have been computed since all indices have been encountered.

## 6 Computational Comparison

In this section we compare the performance of our modifications against the original SMO algorithm and the SMO algorithm that uses duality gap ideas for optimizing  $\beta$ . We implemented all these methods in Fortran and ran them using `f77` on a 200 MHz Pentium machine. The value  $\tau = 0.001$  was used for all experiments. Although we have tested on a number of problems, here we report the performance on only three problems: Wisconsin Breast Cancer data (Bennett & Mangasarian, 1992), UCI Adult data (Platt, 1998), and Web page classification data (Joachims, 1998; Platt, 1998). In the case of the Adult and Web data sets, the inputs are represented in a special binary format, as used by Platt in his testing of SMO. To study scaling properties as training data grow, Platt did staged experiments on the Adult and Web data. We have used only the data from the seventh stage. The gaussian kernel,  $k(x_i, x_j) = \exp(-0.5\|x_i - x_j\|^2/\sigma^2)$ , was used in all experiments. The  $\sigma^2$  values employed, together with  $n$ , the dimension of the input, and  $m$ , the number of training points, are given in Table 1.

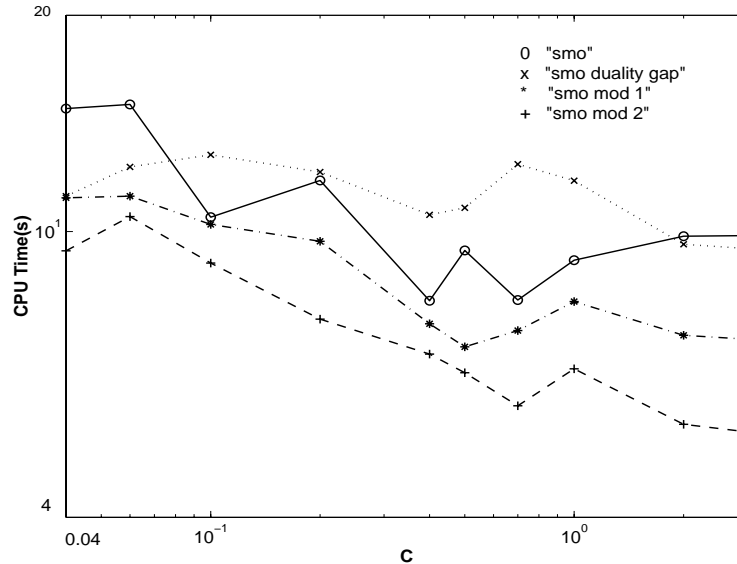


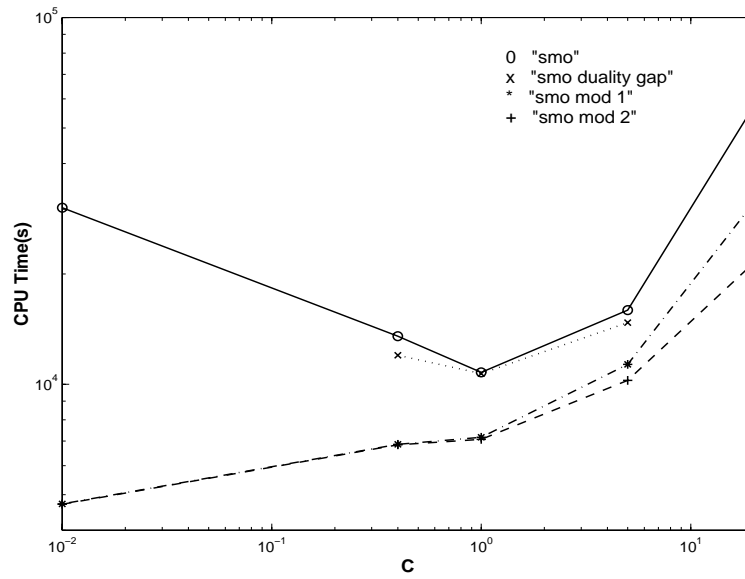
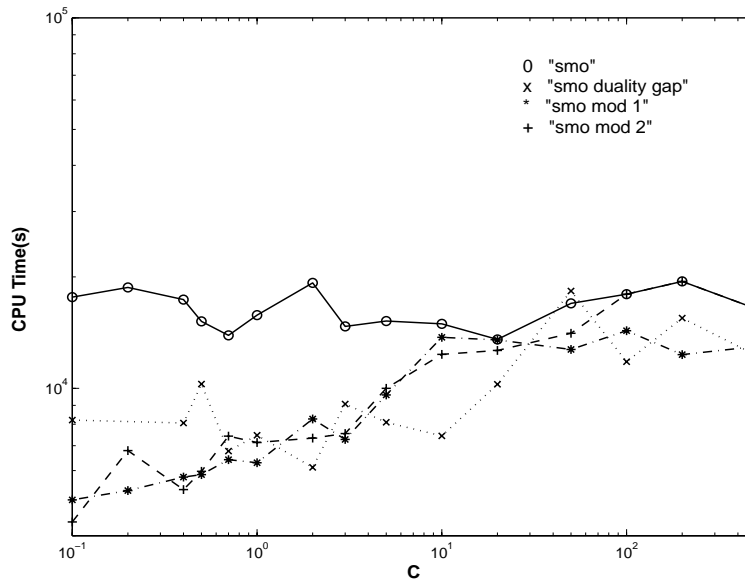
Figure 1: Wisconsin Breast Cancer data. CPU time (in seconds) is shown as a function of  $C$ .

When a particular method is used for SVM design, the value of  $C$  is usually unknown. It has to be chosen by trying a number of values and using a validation set. Therefore, good performance of a method over a range of  $C$  values is important. So for each problem, we tested the algorithms over an appropriate range of  $C$  values. Figures 1 through 3 give the computational costs for the three problems. For the Adult data set, the SMO algorithm based on duality gap ideas had to be terminated for extreme  $C$  values since excessive computing times were needed. In the case of Web data, this method failed for  $C = 0.2$ .

It is clear that the modifications outperform the original SMO algorithm. In most situations, the improvement in efficiency is very good. Between the two modifications, the second one fares better overall. The SMO algorithm based on duality gap did not fare well overall. Although in some cases it did well, in other cases it became slower than the original SMO, and it failed in some cases.

## 7 Conclusion

In this article we have pointed out an important source of inefficiency in Platt's SMO algorithm that is caused by the operation with a single threshold value. We have suggested two modifications of the SMO algorithm that overcome the problem by efficiently maintaining and updating two thresh-

Figure 2: Adult-7 data. CPU time (in seconds) shown as a function of  $C$ .Figure 3: Web-7 data. CPU time (in seconds) shown as a function of  $C$ .

old parameters. Our computational experiments show that these modifications speed up the SMO algorithm considerably in many situations. Platt has already established that the SMO algorithm is one of the fastest algorithms for SVM design. The modified versions of SMO presented here enhance the value of the SMO algorithm even further. It should also be mentioned that heuristics such as shrinking and kernel caching that are effectively used in Joachims (1998) can be employed to speed up our SMO modifications even more. The ideas mentioned in this article for SVM classification can also be extended to the SMO regression algorithm (Smola, 1998). We have reported those results in Shevade et al. (1999).

## References

- Bennett, R., & Mangasarian, O. L. (1992). Robust linear programming discrimination of two linearly inseparable sets. *Optimization Methods and Software*, 1, 23–34.
- Burges, C. J. C. (1998). A tutorial on support vector machines for pattern recognition. *Data Mining and Knowledge Discovery*, 3(2).
- Friess, T. T. (1998). *Support vector networks: The kernel Adatron with bias and soft-margin* (Tech. Rep.). Sheffield, England: University of Sheffield, Department of Automatic Control and Systems Engineering.
- Joachims, T. (1998). Making large-scale support vector machine learning practical. In B. Schölkopf, C. Burges, & A. Smola (Eds.), *Advances in kernel methods: Support vector machines*. Cambridge, MA: MIT Press.
- Keerthi, S. S., Shevade, S. K., Bhattacharyya, C., & Murthy, K. R. K. (1999a). *A fast iterative nearest point algorithm for support vector machine classifier design* (Tech. Rep. No. TR-ISL-99-03). Bangalore, India: Intelligent Systems Lab, Department of Computer Science and Automation, Indian Institute of Science. Available online at: <http://guppy.mpe.nus.edu.sg/~mpessk>.
- Keerthi, S. S., Shevade, S. K., Bhattacharyya, C., & Murthy, K. R. K. (1999b). *Improvements to Platt's SMO algorithm for SVM classifier design* (Tech. Rep. No. CD-99-14). Singapore: Control Division, Department of Mechanical and Production Engineering, National University of Singapore. Available online at: <http://guppy.mpe.nus.edu.sg/~mpessk>.
- Mangasarian, O. L., & Musicant, D. R. (1998). *Successive overrelaxation for support vector machines* (Tech. Rep.). Madison, WI: Computer Sciences Department, University of Wisconsin.
- Platt, J. C. (1998). Fast training of support vector machines using sequential minimal optimization. In B. Schölkopf, C. Burges, & A. Smola (Eds.), *Advances in kernel methods: Support vector machines*. Cambridge, MA: MIT Press.
- Platt, J. C. (1999). Using sparseness and analytic QP to speed training of support vector machines. In M. S. Kearns, S. A. Solla and D. A. Cohn (Eds.), *Advances in neural information processing systems*, 11. Cambridge, MA: MIT Press.
- Shevade, S. K., Keerthi, S. S., Bhattacharyya, C., & Murthy, K. R. K. (1999). *Improvements to the SMO algorithm for SVM regression* (Tech. Rep. No. CD-99-16). Singapore: Control Division, Department of Mechanical and Pro-

duction Engineering, National University of Singapore. Available online at:  
<http://guppy.mpe.nus.edu.sg/~mpessk>.

Smola, A. J. (1998). *Learning with kernels*. Unpublished doctoral dissertation, Technische Universität Berlin.

Vapnik, V. (1995). *The nature of statistical learning theory*. Berlin: Springer-Verlag.

---

Received August 19, 1999; accepted May 17, 2000.