

# A Fast Iterative Nearest Point Algorithm for Support Vector Machine Classifier Design

S.S. Keerthi \*

Dept. of Mechanical and Production Engineering  
National University of Singapore  
Singapore 119260

S.K. Shevade,<sup>†</sup> C. Bhattacharyya<sup>‡</sup> and K.R.K. Murthy<sup>§</sup>

Dept. of Computer Science and Automation  
Indian Institute of Science  
Bangalore, India 560012

April 22, 1999

## Abstract

In this paper we give a new, fast iterative algorithm for support vector machine (SVM) classifier design. The basic problem treated is one that does not allow classification violations. The problem is converted to a problem of computing the nearest point between two convex polytopes. The suitability of two classical nearest point algorithms, due to Gilbert, and Mitchell, Dem'yanov and Malozemov, is studied. Ideas from both these algorithms are combined and modified to derive our fast algorithm. For problems which require classification violations to be allowed, the violations are quadratically penalized and an idea due to Friess is used to convert it to a problem in which there are no classification violations. Comparative computational evaluation of our algorithm against powerful SVM methods such as Platt's Sequential Minimal Optimization shows that our algorithm is very competitive.

## 1. Introduction.

The last few years have seen the rise of Support Vector Machines (SVMs) [25] as powerful tools

---

\* mpessk@guppy.mpe.nus.edu.sg

<sup>†</sup> shirish@csa.iisc.ernet.in

<sup>‡</sup> cbchiru@csa.iisc.ernet.in

<sup>§</sup> murthy@csa.iisc.ernet.in

for solving classification and regression problems[5]. Recently, fast iterative algorithms that are also easy to implement have been suggested[20, 12, 19, 8, 23]; Platt's SMO algorithm [20] is an important example. Such algorithms are bound to widely increase the popularity of SVMs among practitioners. This paper makes another contribution in this direction. Transforming a particular SVM classification problem formulation into a problem of computing the nearest point between two convex polytopes in the hidden feature space, we give a fast iterative nearest point algorithm for SVM classifier design that is competitive with the SMO algorithm. Though not as easy as SMO, our algorithm also is quite straightforward to implement. A pseudocode can be found in [14]. Using this pseudocode an actual running code can be developed within a few hours. A Fortran code implemented by us can be obtained from us for any non-commercial purpose.

The basic problem addressed in this paper is the two category classification problem. Throughout this paper we will use  $x$  to denote the input vector of the support vector machine and  $z$  to denote the feature space vector which is related to  $x$  by a transformation,  $z = \phi(x)$ . As in all 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.

The simplest SVM formulation is one which does not allow classification violations. Let  $\{x_i\}_{i=1}^m$  be a training set of input vectors. Let:  $I$  = the index set for class 1, and  $J$  = the index set for class 2. We assume:  $I \neq \emptyset$ ,  $J \neq \emptyset$ ,  $I \cup J = \{1, \dots, m\}$  and  $I \cap J = \emptyset$ . Let us define:  $z_i = \phi(x_i)$ . The SVM design problem without violations is:

$$\begin{aligned} \min \quad & \frac{1}{2} \|w\|^2 \\ \text{s.t.} \quad & w \cdot z_i + b \geq 1 \quad \forall i \in I; \quad w \cdot z_j + b \leq -1 \quad \forall j \in J \end{aligned} \tag{SVM - NV}$$

Let us make the following assumption:

**Assumption A1.** There exists a  $(w, b)$  pair for which the constraints of SVM-NV are satisfied.

If this assumption holds then SVM-NV has an optimal solution, which turns out to be unique. Let  $H_+ = \{z : w \cdot z + b = 1\}$  and  $H_- = \{z : w \cdot z + b = -1\}$ , the bounding hyperplanes separating the two classes.  $M$ , the margin between them is given by  $M = 2/\|w\|$ . Thus SVM-NV consists of finding the pair of parallel hyperplanes that has the maximum margin among all pairs that separate the two classes.

To deal with data which are linearly inseparable in the  $z$ -space, and also for the purpose of improving generalization, there is a need to have a problem formulation in which classification violations are allowed. The popular approach for doing this is to allow violations in the satisfaction of the constraints in SVM-NV, and penalize such violations *linearly* in the objective function:

$$\begin{aligned} \min \quad & \frac{1}{2} \|w\|^2 + C \sum_k \xi_k \\ \text{s.t.} \quad & w \cdot z_i + b \geq 1 - \xi_i \quad \forall i \in I; \quad w \cdot z_j + b \leq -1 + \xi_j \quad \forall j \in J \\ & \xi_k \geq 0 \quad \forall k \in I \cup J \end{aligned} \tag{SVM - VL}$$

(Throughout, we will use  $k$  and/or  $l$  whenever the indices run over all elements of  $I \cup J$ .) Here  $C$  is a positive, inverse regularization constant that is chosen to give the correct relative weighting between margin maximization and classification violation. Using Wolfe duality theory[5,7] SVM-VL can be transformed to the following equivalent dual problem:

$$\begin{aligned} \max \quad & \sum_k \alpha_k - \frac{1}{2} \sum_k \sum_l \alpha_k \alpha_l y_k y_l z_k \cdot z_l \\ \text{s.t.} \quad & 0 \leq \alpha_k \leq C \quad \forall k; \quad \sum_k \alpha_k y_k = 0, \end{aligned} \quad (\text{SVM - VL - DUAL})$$

where  $y_i = 1, i \in I$  and  $y_j = -1, j \in J$ . It is computationally easy to handle this problem since it is directly based on  $z_i \cdot z_j$  (kernel) calculations. Platt's SMO algorithm, as well as others are ways of solving SVM-VL-DUAL. SMO is particularly very simple and, at the same time, impressively fast. In very simple terms, it is an iterative scheme in which only two (appropriately chosen)  $\alpha_i$  variables are adjusted at any given time so as to improve the value of the objective function. The need for adjusting at least two variables at a time is caused by the presence of the equality constraint.

**Remark 1.** It is useful to point out that the Wolfe dual of SVM-NV is same as SVM-VL-DUAL, with  $C$  set to  $\infty$ . Therefore, any algorithm designed for solving SVM-VL-DUAL can be easily used to solve the dual of SVM-NV, and thereby solve SVM-NV.  $\square$

In a recent paper Friess[9] has suggested the use of a sum of squared violations in the cost function:

$$\begin{aligned} \min \quad & \frac{1}{2} \|w\|^2 + \frac{\tilde{C}}{2} \sum_k \xi_k^2 \\ \text{s.t.} \quad & w \cdot z_i + b \geq 1 - \xi_i \quad \forall i \in I; \quad w \cdot z_i + b \leq -1 + \xi_i \quad \forall i \in J \end{aligned} \quad (\text{SVM - VQ})$$

Preliminary experiments by Friess have shown this formulation to be promising. Unlike SVM-VL, here there is no need to include non-negativity constraints on  $\xi_k$  for the following reason. Suppose, at the optimal solution of SVM-VQ,  $\xi_k$  is negative for some  $k$ . Then, by resetting  $\xi_k = 0$ , we can remain feasible and also strictly decrease the cost. Thus, negative values of  $\xi_k$  cannot occur at the optimal solution.

**Remark 2.** As pointed out by Friess[9], a very nice property of SVM-VQ is that by doing a simple transformation it can be converted into an instance of SVM-NV. Let  $e_k$  denote the  $m$  dimensional vector in which the  $k$ -th component is 1 and all other components are 0. Define:

$$\tilde{w} = \begin{pmatrix} w \\ \sqrt{\tilde{C}} \xi \end{pmatrix}; \quad b = \tilde{b}; \quad \tilde{z}_i = \begin{pmatrix} z_i \\ \frac{1}{\sqrt{\tilde{C}}} e_i \end{pmatrix}, i \in I; \quad \tilde{z}_j = \begin{pmatrix} z_j \\ -\frac{1}{\sqrt{\tilde{C}}} e_j \end{pmatrix}, j \in J. \quad (1)$$

Then it is easy to see that SVM-V transforms to an instance of SVM-NV. (Use  $\tilde{w}$ ,  $\tilde{b}$ ,  $\tilde{z}$  instead of  $w$ ,  $b$ ,  $z$ .) Note that, because of the presence of  $\xi_k$  variables, SVM-VQ's feasible space is always non-empty and so the resulting SVM-NV is automatically feasible. Also note that, if  $K$  denotes the Kernel function in the SVM-VQ problem, then  $\tilde{K}$ , the Kernel function for the transformed

SVM-NV problem is given by

$$\tilde{K}(x_k, x_l) = K(x_k, x_l) + \frac{1}{C} \delta_{kl}$$

where  $\delta_{kl}$  is 1 if  $k = l$  and 0 otherwise. Thus, for any pair of training vectors,  $(x_k, x_l)$ , the modified kernel function is easily computed.  $\square$

Our main aim in this paper is to give a fast algorithm for solving SVM-NV. The main idea consists of transforming SVM-NV into a problem of computing the nearest point between two convex polytopes and then using a carefully chosen nearest point algorithm to solve it. Because of Remark 2, our algorithm can also be easily used to solve SVM-VQ. By Remark 1, algorithms such as SMO can also be used to solve SVM-VQ. In empirical testing however, we have found that (details are given in section 8) our algorithm is more efficient for solving SVM-VQ, than SMO used in this way. Even when the “typical” computational cost of solving SVM-VL by SMO is compared with that of solving SVM-VQ by our algorithm, we find that our algorithm is competitive.

Recently Friess et.al., [8] (deriving inspiration from the Adatron algorithm given by Anlauf and Biehl[1] for designing Hopfield nets) and Mangasarian and Musicant[19] (using a successive overrelaxation idea) independently suggested the inclusion of the extra term,  $b^2/2$  in the objective functions of the various formulations. This is done with computational simplicity in mind. When  $b^2/2$  is added to the objective functions of the primal SVM problems, i.e., SVM-NV, SVM-VL, and SVM-VQ, it turns out that the only equality constraint in the dual problems, i.e.,  $\sum_i \alpha_i y_i = 0$ , gets eliminated. Therefore, it is easy to give an iterative algorithm for improving the dual objective function simply by adjusting a single  $\alpha_i$  at a time, as opposed to the adjustment of two such variables required by SMO. Friess et.al.[8] applied the Adatron algorithm to solve SVM-NV, Friess[9] applied the same to SVM-VQ, while Mangasarian and Musicant[19] applied the successive overrelaxation scheme to solve SVM-VL. The basic algorithmic ideas used by them is as follows: choose one  $\alpha_i$ , determine a unconstrained step size for  $\alpha_i$  as if there are no bounds on  $\alpha_i$ , and then clip the step size so that the updated  $\alpha_i$  satisfies all its bounds.

If the term  $b^2/2$  is included in the objective functions of SVM-VQ and SVM-NV, then these problems can be easily transformed to a problem of computing the nearest point of a single convex polytope from the origin, a problem that is much simpler than finding the nearest distance between two convex polytopes. Our nearest point algorithm mentioned earlier simplifies considerably for this simpler problem. Our testing shows that these simplified algorithms do not perform as well as algorithms which solve problems without the  $b^2/2$  term.

The paper is organized as follows. In section 2 we reformulate SVM-NV as a problem of computing the nearest point between two convex polytopes. In section 3 we discuss optimality criteria for this nearest point problem. Section 4 discusses important practical issues concerning approximate stopping of dual solutions. It points out the danger of simply stopping the dual solution using dual objective function accuracy, and derives a simple check for stopping nearest point

algorithms so as to get guaranteed accuracy for the solution of SVM-NV. The main algorithm of the paper is derived in sections 5 and 6. Two classical algorithms for doing nearest point solution, due to Gilbert[10] and Mitchell et.al.[18], are combined and modified to derive our fast algorithm. Section 7 concerns the simplification of this algorithm to treat the case corresponding to the inclusion of  $b^2/2$  in the objective functions of SVM-VQ and SVM-NV. A detailed comparative computational testing of existing algorithms against our algorithms is taken up in section 8. We conclude with some closing remarks in section 9.

## 2. Reformulation of SVM-NV as a Nearest Point Problem.

Given a set  $S$  we use  $\text{co}S$  to denote the convex hull of  $S$ , i.e.,  $\text{co}S$  is the set of all convex combinations of elements of  $S$ :

$$\text{co}S = \left\{ \sum_{k=1}^l \beta_k s_k : s_k \in S, \beta_k \geq 0, \sum_{k=1}^l \beta_k = 1 \right\}$$

Let  $U = \text{co}\{z_i : i \in I\}$  and  $V = \text{co}\{z_i : i \in J\}$  where  $\{z_k\}$ ,  $I$  and  $J$  are as in the definition of SVM-NV. Since  $I$  and  $J$  are finite sets,  $U$  and  $V$  are convex polytopes. Consider the following generic problem of computing the minimum distance between  $U$  and  $V$ :

$$\min \|u - v\| \quad \text{s.t. } u \in U, v \in V \quad (\text{NPP})$$

It can be easily noted that the solution of NPP may not be unique. We can rewrite the constraints of NPP algebraically as:

$$\begin{aligned} u &= \sum_{i \in I} \beta_i z_i; & \beta_i &\geq 0, i \in I; & \sum_{i \in I} \beta_i &= 1 \\ v &= \sum_{j \in J} \beta_j z_j; & \beta_j &\geq 0, j \in J; & \sum_{j \in J} \beta_j &= 1 \end{aligned} \quad (2)$$

The equivalence of the problems SVM-NV and NPP can be easily understood by studying the geometry shown in Figure 1. Assumption A1 is equivalent to assuming that  $U$  and  $V$  are non-intersecting. Thus A1 implies that the optimal cost of NPP is positive. If  $(w^*, b^*)$  denotes the solution of SVM-NV and  $(u^*, v^*)$  denotes a solution of NPP, then by using the facts that maximum margin  $= 2/\|w^*\| = \|u^* - v^*\|$  and  $w^* = \delta(u^* - v^*)$  for some  $\delta$ , we can easily derive the following relationship between the solutions of SVM-NV and NPP.

$$w^* = \frac{2}{\|u^* - v^*\|^2} (u^* - v^*); \quad b^* = \frac{\|v^*\|^2 - \|u^*\|^2}{\|u^* - v^*\|^2} \quad (3)$$

The following theorem states this relationship formally.

**Theorem 1.**  $(w^*, b^*)$  solves SVM-NV if and only if there exist  $u^* \in U$  and  $v^* \in V$  such that  $(u^*, v^*)$  solves NPP and (3) holds.  $\square$

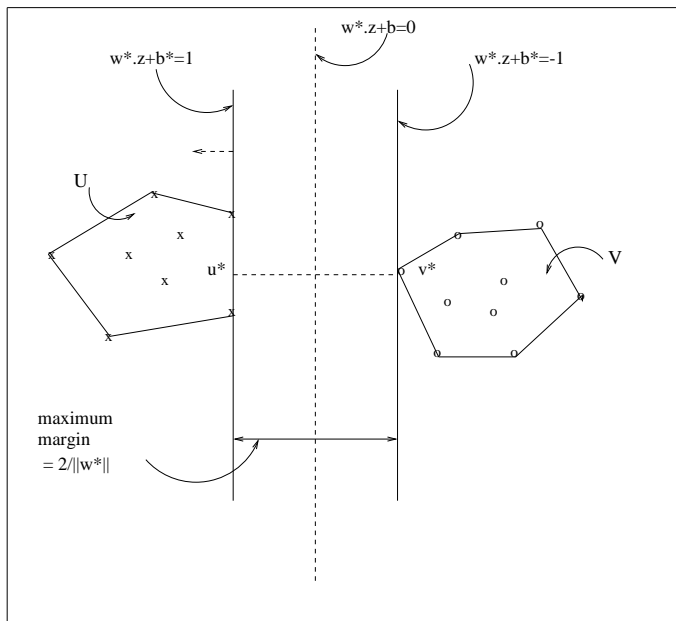


Figure 1: Among all pairs of parallel hyperplanes that separate the two classes, the pair with the largest margin is the one which has  $(u^*, v^*)$  as the normal direction, where  $(u^*, v^*)$  is a pair of closest points of  $U$  and  $V$ . Note that  $w^* = \delta(u^* - v^*)$  for  $\delta$  chosen such that  $\|u^* - v^*\| = 2/\|w^*\|$ .

A direct, geometrically intuitive proof is given by Sancheti and Keerthi[22] with reference to a geometrical problem in Robotics. Later, Bennett[3] proved a somewhat close result in the context of learning algorithms. Here we only give a discussion that follows the traditional Wolfe-Dual approach employed in the SVM literature. The main reason for doing this is to show the relationships of NPP and the variables in it with the Wolfe-Dual of SVM-NV and the variables there.

Using Wolfe duality theory[5,7] we first transform SVM-NV to the following equivalent dual problem:

$$\begin{aligned} \max \quad & \sum_k \alpha_k - \frac{1}{2} \sum_k \sum_l \alpha_k \alpha_l y_k y_l z_k \cdot z_l \\ \text{s.t.} \quad & \alpha_k \geq 0 \quad \forall k; \quad \sum_k \alpha_k y_k = 0, \end{aligned} \quad (\text{SVM - NV - DUAL})$$

where, as before,  $y_i = 1, i \in I$  and  $y_j = -1, j \in J$ . Since  $\sum_k \alpha_k y_k = 0$  implies that  $\sum_{i \in I} \alpha_i = \sum_{i \in J} \alpha_i$ , we can introduce a new variable,  $\lambda$  and rewrite  $\sum_k \alpha_k y_k = 0$  as two constraints:

$$\sum_{i \in I} \alpha_i = \lambda, \quad \sum_{i \in J} \alpha_i = \lambda$$

If we also define

$$\beta_k = \frac{\alpha_k}{\lambda} \quad \forall k \quad (4)$$

then SVM-NV-DUAL can be rewritten as

$$\begin{aligned} \max \quad & 2\lambda - \frac{\lambda^2}{2} \sum_k \sum_l \beta_k \beta_l y_k y_l z_k \cdot z_l \\ \text{s.t.} \quad & \beta_k \geq 0 \quad \forall k; \quad \sum_{i \in I} \beta_i = 1; \quad \sum_{j \in J} \beta_j = 1. \end{aligned} \quad (5)$$

In this formulation it is convenient to first optimize  $\lambda$  keeping  $\beta$  constant, and then optimize  $\beta$  on the outer loop. Optimizing with respect to  $\lambda$  yields

$$\lambda^* = \frac{2}{\sum_k \sum_l \beta_k \beta_l y_k y_l z_k \cdot z_l} \quad (6)$$

Substituting this in (5) and simplifying, we see that (5) becomes equivalent to:

$$\begin{aligned} \max \quad & 2 / (\sum_k \sum_l \beta_k \beta_l y_k y_l z_k \cdot z_l) \\ \text{s.t.} \quad & \beta_k \geq 0 \quad \forall k; \quad \sum_{i \in I} \beta_i = 1; \quad \sum_{j \in J} \beta_j = 1. \end{aligned} \quad (7)$$

This problem is equivalent to the problem:

$$\begin{aligned} \min \quad & \frac{1}{2} \sum_k \sum_l \beta_k \beta_l y_k y_l z_k \cdot z_l \\ \text{s.t.} \quad & \beta_k \geq 0 \quad \forall k; \quad \sum_{i \in I} \beta_i = 1; \quad \sum_{j \in J} \beta_j = 1. \end{aligned} \quad (8)$$

If we define a matrix  $P$  whose columns are  $y_1 z_1, \dots, y_m z_m$  then it is easy to note that

$$\sum_k \sum_l \beta_k \beta_l y_k y_l z_k \cdot z_l = \|P\beta\|^2$$

If  $\beta$  satisfies the constraints of (8), then  $P\beta = u - v$  where  $u \in U$  and  $v \in V$ . Thus (8) is equivalent to NPP.

Wolfe duality theory[5] actually yields

$$w^* = \sum_k \alpha_k^* y_k z_k$$

Using this, (4) and (6) we can immediately verify the expression for  $w^*$  in (3). The expression for  $b^*$  can be easily understood from geometry.

**Remark 3.** The above discussion also points out an important fact: there is a simple redundancy in the SVM-NV-DUAL formulation which gets removed in the NPP formulation. Note that NPP has two equality constraints and SVM-NV-DUAL has only one, while both have the same number of variables. Therefore, when quadratic programming algorithms are applied to SVM-NV-DUAL and NPP separately they work quite differently, even when started from the same “equivalent” starting points.  $\square$

### 3. Optimality Criteria for NPP.

First let us give some definitions concerning support properties of a convex polytope. For a given compact set  $P$ , let us define the *support function*,  $h_P : R^n \rightarrow R$ , by

$$h_P(\eta) = \max\{\eta \cdot z : z \in P\} \quad (9)$$

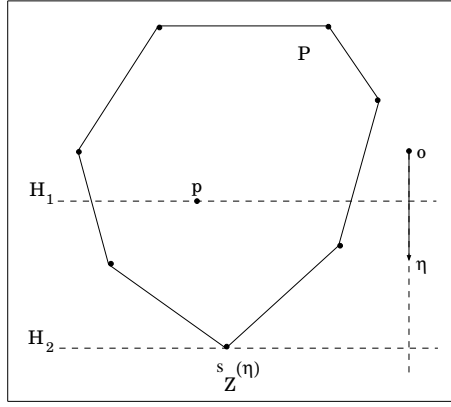


Figure 2: Definition of support properties.  $s_Z(\eta)$  is the extreme vertex of  $P$  in the direction  $\eta$ . The distance between the hyperplanes,  $H_1$  and  $H_2$  is equal to  $g_P(\eta, p)/\|\eta\|$ .

See Figure 2. We use  $s_P(\eta)$  to denote any one solution of (9), i.e.,  $s_P(\eta)$  satisfies

$$h_P(\eta) = s_P(\eta) \cdot \eta \text{ and } s_P(\eta) \in P \quad (10)$$

Now consider the case where  $P$  is a convex polytope:  $Z = \{\hat{z}_1, \dots, \hat{z}_r\}$  and  $P = \text{co}Z$ . It is well-known[16] that the maximum of a linear function over a convex polytope is attained by an extreme point. This means that  $h_P = h_Z$  and  $s_P = s_Z$ . Therefore  $h_P$  and  $s_P$  can be determined by a simple enumeration of inner products:

$$\begin{aligned} h_P(\eta) &= h_Z(\eta) = \max\{\eta \cdot \hat{z}_k : k = 1, \dots, r\} \\ s_P(\eta) &= s_Z(\eta) = \hat{z}_l \text{ where } \eta \cdot \hat{z}_l = h_Z(\eta) \end{aligned} \quad (11)$$

Thus (11) provides a simple procedure for evaluating the support properties of the convex polytope,  $P = \text{co}Z$ . Let us also define the function,  $g_P : R^n \times P \rightarrow R$  by

$$g_P(\eta, p) = h_P(\eta) - \eta \cdot p$$

By (9) it follows that

$$g_P(\eta, p) \geq 0 \quad \forall \eta \in R^n, p \in P \quad (12)$$

We now adapt these general definitions to derive an optimality criterion for NPP. A simple geometrical analysis shows that a pair,  $(u, v) \in U \times V$  solves NPP if and only if

$$h_U(-z) = -z \cdot u \quad \text{and} \quad h_V(z) = z \cdot v$$

where  $z = u - v$ ; in other words,  $(u, v)$  solves NPP if and only if  $u = s_U(-z)$  and  $v = s_V(z)$ . Equivalently,  $(u, v)$  is optimal if and only if

$$g_U(v - u, u) = 0 \quad \text{and} \quad g_V(u - v, v) = 0 \quad (13)$$

Let us define the function,  $g : U \times V \rightarrow R$  by

$$g(u, v) = g_U(v - u, u) + g_V(u - v, v) \quad (14)$$

Since  $g_U$  and  $g_V$  are nonnegative functions, it also follows that  $(u, v)$  is optimal if and only if  $g(u, v) = 0$ . The following theorem states the above results (and related ones) formally.

**Theorem 2.** Suppose  $u \in U$ ,  $v \in V$  and  $z = u - v$ . Then the following hold. (a)  $g(u, v) \geq 0$ . (b) if  $\bar{u} \in U$  is a point that satisfies  $z \cdot \bar{u} < z \cdot u$ , then there is a point  $\tilde{u}$  on the line segment,  $\text{co}\{u, \bar{u}\}$  such that  $\|\tilde{u} - v\| < \|u - v\|$ . (c) if  $\bar{v} \in V$  is a point that satisfies  $z \cdot \bar{v} < z \cdot v$ , then there is a point  $\tilde{v}$  on the line segment,  $\text{co}\{v, \bar{v}\}$  such that  $\|u - \tilde{v}\| < \|u - v\|$ . (d)  $(u, v)$  solves NPP if and only if  $g(u, v) = 0$ .

**Proof.** (a) This follows from (12) and (14).

(b) Define a real variable  $t$  and a function,  $\phi(t) = \|u + t(\bar{u} - u) - v\|^2$ .  $\phi$  describes the variation of  $\|\hat{u} - v\|^2$  as a generic point  $\hat{u}$  varies from  $u$  to  $\bar{u}$  over the line segment joining them. Since  $\phi'(0) = 2z \cdot (\bar{u} - u) < 0$ , the result follows.

(c) The proof is along similar lines as that of (b).

(d) First let  $g(u, v) = 0$ . Let  $\hat{u} \in U$ ,  $\hat{v} \in V$ . Since  $g_U(-z, u) = 0$  and  $g_V(z, v) = 0$ , we have  $z \cdot u \leq z \cdot \hat{u}$  and  $z \cdot v \geq z \cdot \hat{v}$ . Subtracting the two inequalities we get  $\|z\|^2 \leq z \cdot \hat{z}$  where  $\hat{z} = \hat{u} - \hat{v}$ . Now

$$\|z\|^2 \leq \|z\|^2 + \|z - \hat{z}\|^2 = \|\hat{z}\|^2 + 2(\|z\|^2 - \hat{z} \cdot z) \leq \|\hat{z}\|^2$$

and hence  $(u, v)$  is optimal for NPP. On the other hand, if  $(u, v)$  is optimal for NPP, we can set  $\bar{u} = s_U(-z)$ ,  $\bar{v} = s_V(z)$  and use results (a) and (b) to show that  $g(u, v) = 0$ .  $\square$

#### 4. Stopping Criteria for NPP Algorithms.

When a numerical algorithm is employed to solve any of the problems mentioned earlier, approximate stopping criteria need to be employed. In SVM design a dual formulation (such as SVM-NV-DUAL and NPP) is solved instead of the primal (such as SVM-NV) because of computational ease. However it has to be kept in mind that it is the primal solution that is of final interest and that care is needed to ensure that the numerical algorithm has reached an approximate primal solution with a guaranteed specified closeness to the optimal primal solution. We will show using an example that simply stopping when the dual approximate solution has reached the dual optimal solution within a good accuracy can be dangerous.

Consider the example shown in Figure 3. Let us take it that Class 2 has a single point, which is the origin,  $o$ . Thus  $v^* = o$ . Let  $u^*$  be on the line segment joining a and b, which are the support vectors from Class 1 (i.e., those vectors in  $\{z_i\}$  which have positive multipliers in a representation of  $u^*$ ). Let us use the notation,  $\overline{xy}$  to denote the distance between the two points,  $x$  and  $y$ . Let

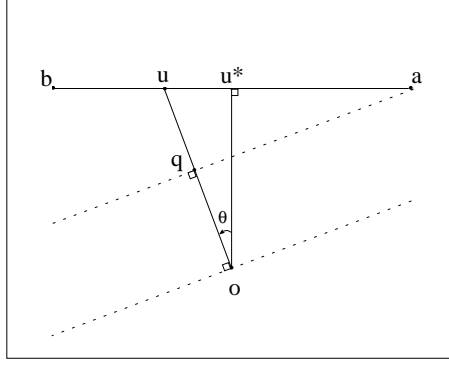


Figure 3: Example to illustrate the danger of stopping with only good dual solution accuracy.

$\rho^* = \overline{ou^*}$ , the optimal distance between the two classes, and,  $d = \overline{oa}$ . Suppose it turns out that  $d$  is very much larger than  $\rho^*$ . (This condition is quite typical of most classification problems.) Suppose we are solving NPP and we stop when we reach an approximate solution,  $u$  that satisfies  $\overline{ou} \leq (1 + \epsilon)\overline{ou^*}$ . Let us assume that  $\epsilon$  is much smaller than 1. If we assume that the support vectors are identified correctly (most algorithms, in fact do this well) then  $u$  lies on the line segment joining the two support vectors,  $a$  and  $b$ . Let  $u$  be as shown in Figure 3.

Consider triangle  $ouu^*$ .

$$\overline{uu^*} = \sqrt{(\overline{ou})^2 - (\overline{ou^*})^2} = \sqrt{(2\epsilon + \epsilon^2)\rho^*} \approx \sqrt{2\epsilon}\rho^*$$

Then  $\sin \theta = \overline{uu^*}/\overline{ou} = \sqrt{2\epsilon + \epsilon^2}/(1 + \epsilon) \approx \sqrt{2\epsilon}$ . Now consider triangle  $ou^*a$ .  $\overline{au^*} = \sqrt{(\overline{ao})^2 - (\overline{ou^*})^2} = \sqrt{d^2 - (\rho^*)^2} \approx d$ . Also,  $\overline{au} = \overline{au^*} + \overline{uu^*} \approx d$ . Let  $q$  be the point of intersection of the line joining  $o$  and  $u$  with the line which has  $\overline{ou^*}$  as the normal direction and passes through  $a$ . Clearly  $\overline{oq}$  is  $M$ , the margin corresponding to the approximate solution,  $u$ . Note that  $M^*$ , the optimal margin for the problem, is equal to  $\rho^*$ . To get an estimate of  $M$  consider triangle  $uqa$ .  $\overline{uq} = \overline{au} \sin \theta \approx d\sqrt{2\epsilon}$ . Then,

$$M = \overline{ou} - \overline{uq} \approx \rho^*\sqrt{1 + \epsilon} - d\sqrt{2\epsilon} \approx (1 + 0.5\epsilon)\rho^* - d\sqrt{2\epsilon}$$

From this we finally get

$$\frac{(M^* - M)}{M^*} = \frac{(\rho^* - M)}{\rho^*} \approx \frac{d}{\rho^*}\sqrt{2\epsilon} - 0.5\epsilon$$

For small  $\epsilon$ , it is the first term on the right hand side that dominates. Thus, when  $d/\rho^*$  is large, the error in the margin can be very large! A small deviation,  $u$  from  $u^*$  thus leads to a large change in the margin,  $M$  from  $M^*$ , simply because of the reason that  $a$  is much farther from  $o$  than  $u^*$ ! What is dangerous is that, even with a small  $\epsilon$ ,  $M$  can end-up being negative and hence, at the end of the algorithm we cannot even generate a feasible  $(w, b)$  pair for SVM-NV using the approximate dual solution. Of course, when  $\epsilon$  is made very very small, these problems disappear. However, a key

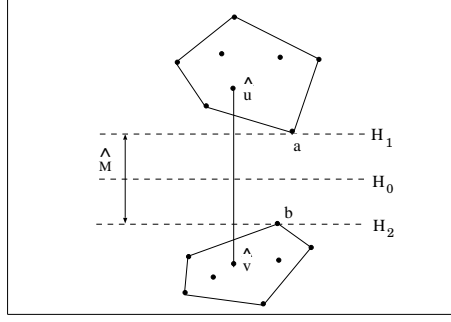


Figure 4: A situation satisfying (15). Here:  $a = s_U(-\hat{z})$ ;  $b = s_V(\hat{z})$ ;  $H_1 : \hat{z} \cdot z = -h_U(-\hat{z})$ ;  $H_2 : \hat{z} \cdot z = h_V(\hat{z})$ ; and  $H_0 : \hat{z} \cdot z = 0.5(h_V(\hat{z}) - h_U(-\hat{z}))$ .

question still remains to be answered: how do we systematically stop the dual solution algorithm in order to produce a feasible primal solution that has a specified closeness to the optimal primal solution? We will next show that the  $g$  function that we introduced in (14) greatly aids in doing such a guaranteed stopping.

Suppose we have an algorithm for NPP that iteratively improves its approximate solution,  $\hat{u} \in U$ ,  $\hat{v} \in V$  in such a way that  $\hat{u} \rightarrow u^*$  and  $\hat{v} \rightarrow v^*$ , where  $(u^*, v^*)$  is a solution of NPP. By part (d) of Theorem 2 and the fact that  $g$  is a continuous function,  $g(\hat{u}, \hat{v}) \rightarrow 0$ . Also,  $\|\hat{z}\| \rightarrow \|z^*\|$  where  $\hat{z} = \hat{u} - \hat{v}$  and  $z^* = u^* - v^*$ . By assumption A1,  $\|\hat{z}\| \geq \|z^*\| > 0$ . Thus we also have  $g(\hat{u}, \hat{v})/\|\hat{z}\|^2 \rightarrow 0$ . Suppose  $\epsilon$  is a specified accuracy parameter satisfying  $0 < \epsilon < 1$ , and that it is desired to stop the algorithm when we have found a  $(\hat{w}, \hat{b})$  pair that is feasible to SVM-NV and  $\|w^*\| \geq (1 - \epsilon)\|\hat{w}\|$  where  $(w^*, b^*)$  is the optimal solution of SVM-NV. We will show that this is possible if we stop the dual algorithm when  $(\hat{u}, \hat{v})$  satisfies

$$g(\hat{u}, \hat{v}) \leq \epsilon \|\hat{z}\|^2 \quad (15)$$

Figure 4 geometrically depicts the situation. Let  $\hat{M}$  be the margin corresponding to the direction,  $\hat{z}$ . Clearly,

$$\hat{M} = \frac{-h_U(-\hat{z}) - h_V(\hat{z})}{\|\hat{z}\|} \quad (16)$$

Note that, since  $g(\hat{u}, \hat{v}) = \|\hat{z}\|^2 + h_U(-\hat{z}) + h_V(\hat{z}) \leq \epsilon \|\hat{z}\|^2$ , we have

$$-h_U(-\hat{z}) - h_V(\hat{z}) \geq (1 - \epsilon) \|\hat{z}\|^2 \quad (17)$$

Thus,  $\hat{M} \geq (1 - \epsilon) \|\hat{z}\| > 0$ . The determination of  $\hat{w}$  that is consistent with the feasibility of SVM-NV can be easily found by setting

$$\hat{w} = \gamma \hat{z} \quad (18)$$

and choosing  $\gamma$  such that  $\hat{M} = 2/\|\hat{w}\|$ . Using (16) we get

$$\gamma = \frac{2}{(-h_U(-\hat{z}) - h_V(\hat{z}))} \quad (19)$$

Since the equation of the central separating hyperplane,  $H_0$  has to be of the form,  $\hat{w} \cdot z + \hat{b} = 0$ , we can also easily get the expression for  $\hat{b}$  as

$$\hat{b} = \frac{h_V(\hat{z}) - h_U(-\hat{z})}{h_V(\hat{z}) + h_U(-\hat{z})} \quad (20)$$

Using (16), (17) and the fact that  $M^* = \|z^*\|$ , we get

$$\frac{\|w^*\|}{\|\hat{w}\|} = \frac{\hat{M}}{M^*} = \frac{-h_U(-\hat{z}) - h_V(\hat{z})}{\|\hat{z}\|\|z^*\|} = \frac{-h_U(-\hat{z}) - h_V(\hat{z})}{\|\hat{z}\|^2} \cdot \frac{\|\hat{z}\|}{\|z^*\|} \geq (1 - \epsilon)$$

To derive a bound for  $\|z^*\|/\|\hat{z}\|$ , note that  $\|z^*\| = M^* \geq \hat{M} = (-h_U(-\hat{z}) - h_V(\hat{z}))/\|\hat{z}\|$ . Then,

$$\frac{\|z^*\|}{\|\hat{z}\|} \geq \frac{-h_U(-\hat{z}) - h_V(\hat{z})}{\|\hat{z}\|^2} \geq (1 - \epsilon)$$

Thus we have proved the following important result.

**Theorem 3.** Let  $(u^*, v^*)$  be a solution of NPP,  $z^* = u^* - v^*$ , and  $(w^*, b^*)$  be the optimal solution of SVM-NV. Let  $0 < \epsilon < 1$ . Suppose  $\hat{u} \in U$ ,  $\hat{v} \in V$ ,  $\hat{z} = \hat{u} - \hat{v}$  and (15) holds. Then  $(\hat{w}, \hat{b})$  as defined by (18)-(20) is feasible for SVM-NV and

$$\min\left(\frac{\|w^*\|}{\|\hat{w}\|}, \frac{\|\hat{M}\|}{\|M^*\|}, \frac{\|z^*\|}{\|\hat{z}\|}\right) \geq (1 - \epsilon)$$

If one is particularly interested in getting a bound on the cost function of SVM-NV then it can be easily obtained:

$$\frac{\|w^*\|^2}{\|\hat{w}\|^2} = \left(\frac{\|w^*\|}{\|\hat{w}\|}\right)^2 \geq (1 - \epsilon)^2 \geq (1 - 2\epsilon)$$

A similar bound can be obtained for  $\|z^*\|^2/\|\hat{z}\|^2$ .

All these discussions point to the important fact that (15) can be used to effectively terminate a numerical algorithm for solving NPP.

## 5. Iterative Algorithms for NPP.

NPP has been well studied in the literature, and a number of good algorithms have been given for it[10, 18, 27, 2, 15]. Best general purpose algorithms for NPP such as Wolfe's algorithm[27] terminate within a finite number of steps; however they require expensive matrix storage and matrix operations in each step that makes them unsuitable for use in large SVM design.<sup>1</sup> Iterative

---

<sup>1</sup>It is interesting to point out that the active set method used by Kaufman[13] to solve SVMs is *identical* to Wolfe's algorithm when both are used to solve SVM-NV.

algorithms that need minimal memory (i.e., memory size needed is linear in  $m$ , the number of training vectors), but which only reach the solution asymptotically as the number of iterations goes to infinity, seem to be better suited for SVM design. In this section we will take up some such algorithms for investigation.

### 5.1. Gilbert's Algorithm.

Gilbert's algorithm[10] was one of the first algorithms suggested for solving NPP. It was originally devised to solve certain optimal control problems. Later its modifications have found good use in pattern recognition and robotics[27,11]. In this section we briefly describe the algorithm and point to how it can be adapted for SVM classifier design.

Let  $\{z_k\}$ ,  $I$ ,  $J$ ,  $U$  and  $V$  be as in the definition of NPP. Let  $Z$  denote the Minkowski set difference of  $U$  and  $V$ [16], i.e.,

$$Z = U \ominus V = \{u - v : u \in U, v \in V\}$$

Clearly, NPP is equivalent to the following minimum norm problem:

$$\min\{\|z\| : z \in Z\} \tag{MNP}$$

Unlike NPP this problem has a unique solution. MNP is very much like NPP and seems like a simpler problem than NPP, but it is only superficially so.  $Z$  is the convex polytope,  $Z = \text{co}W$ , where

$$W = \{z_i - z_j : i \in I, j \in J\} \tag{21}$$

a set with  $m_1 m_2$  points ( $m_1 = |I|$ ,  $m_2 = |J|$ ). To see this, take two general points,  $u \in U$  and  $v \in V$  with the representation,

$$\begin{aligned} u &= \sum_{i \in I} \beta_i z_i, \quad \sum_{i \in I} \beta_i = 1, \quad \beta_i \geq 0 \quad \forall i \in I, \\ v &= \sum_{j \in J} \beta_j z_j, \quad \sum_{j \in J} \beta_j = 1, \quad \beta_j \geq 0 \quad \forall j \in J \end{aligned}$$

Then write  $z = u - v$  as

$$\begin{aligned} z &= \sum_{i \in I} \beta_i z_i - \sum_{j \in J} \beta_j z_j \\ &= (\sum_{j \in J} \beta_j) \sum_{i \in I} \beta_i z_i - (\sum_{i \in I} \beta_i) \sum_{j \in J} \beta_j z_j \\ &= \sum_{i \in I} \sum_{j \in J} \beta_i \beta_j (z_i - z_j) \end{aligned}$$

and note the fact that  $\sum_{i \in I} \sum_{j \in J} \beta_i \beta_j = (\sum_{i \in I} \beta_i)(\sum_{j \in J} \beta_j) = 1$ . In general it is possible for  $Z$  to have  $O(m_1 m_2)$  vertices. Since  $m_1 m_2$  can become very large, it is definitely not a good idea to form  $W$  and then define  $Z$  using it.

Gilbert's algorithm is a method of solving MNP that does not require the explicit formation of  $W$ . Its steps require only the evaluations of the support properties,  $h_Z$  and  $s_Z$ . Fortunately these functions can be computed efficiently:

$$\begin{aligned} h_Z(\eta) = h_W(\eta) &= \max_{i \in I, j \in J} \eta \cdot (z_i - z_j) \\ &= \max_{i \in I} (\eta \cdot z_i) + \max_{j \in J} (-\eta \cdot z_j) = h_U(\eta) + h_V(-\eta) \end{aligned} \quad (22)$$

$$s_Z(\eta) = s_U(\eta) - s_V(-\eta) \quad (23)$$

Thus, for a given  $\eta$  the computation of  $h_Z(\eta)$  and  $s_Z(\eta)$  requires only  $O(m_1 + m_2)$  time.

An optimality criterion for MNP can be easily derived as in section 3. This is stated in the following theorem. We will omit its proof since it is very much along the lines of proof of Theorem 2. The  $g$  function of (12), as applied to  $Z$  plays a key role.

**Theorem 4.** Suppose  $z \in Z$ . Then the following hold. (a)  $g_Z(-z, z) \geq 0$ . (b) If  $\bar{z}$  is any point in  $Z$  such that  $\|z\|^2 - z \cdot \bar{z} > 0$ , there is a point  $\tilde{z}$  in the line segment,  $\text{co}\{z, \bar{z}\}$  satisfying  $\|\tilde{z}\| < \|z\|$ . (c)  $z$  solves MNP if and only if  $g_Z(-z, z) = 0$ .  $\square$

Gilbert's algorithm is based on the results in the above theorem. Suppose we start with some  $z \in Z$ . If  $g_Z(-z, z) = 0$  then  $z = z^*$ , the solution of MNP. On the other hand, if  $g_Z(-z, z) > 0$ , then  $\bar{z} = s_Z(-z)$  satisfies  $\|z\|^2 - z \cdot \bar{z} > 0$ . By part (b) of the theorem, then, searching on the line segment connecting  $z$  to  $\bar{z}$  yields a point whose norm is smaller than  $\|z\|$ . These observations yield Gilbert's algorithm.

#### Gilbert's Algorithm for solving MNP.

0. Choose  $z \in Z$ .
1. Compute  $h_Z(-z)$  and  $g_Z(-z, z)$ . If  $g_Z(-z, z) = 0$  stop with  $z^* = z$ ; else, set  $\bar{z} = s_Z(-z)$ .
2. Compute  $\tilde{z}$ , the point on the line segment joining  $z$  and  $\bar{z}$  which has least norm, set  $z = \tilde{z}$  and go back to step 1.  $\square$

An efficient algorithm for computing the point of least norm on the line segment joining two points is given in the appendix. Figure 5 shows a few iterations of Gilbert's algorithm on a two dimensional example. Gilbert showed that, if the algorithm does not stop with  $z^*$  at step 1 within a finite number of iterations, then  $z \rightarrow z^*$  asymptotically, as the number of iterations goes to infinity. He also showed that the algorithm has a linear rate of convergence.

We now discuss how Gilbert's algorithm can be adapted to solve an NPP that arises from an SVM classifier design problem. Step 0 can be easily implemented by choosing any  $i \in I, j \in J$  and setting  $z = z_i - z_j$ . In a typical usage of Gilbert's algorithm,  $z$  is actually maintained in memory. This is possible for linear SVM design, where the  $x$  space and  $z$  space coincide. However, in a nonlinear SVM design problem that is based on kernel functions,  $z$  is located either in a very large dimensional space or in an infinite dimensional space. In such a case, instead of maintaining  $z$  one

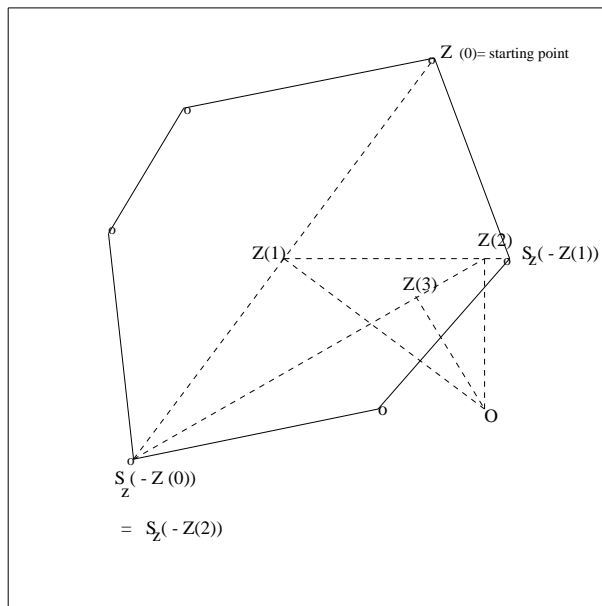


Figure 5: A few iterations of Gilbert's algorithm on a two dimensional example.

has to maintain a convex hull representation of  $z$ :

$$z = \sum_t \gamma_t w_t, \quad \sum_t \gamma_t = 1, \quad \gamma_t \geq 0 \forall t$$

where  $\{w_t\} \subset W$  (see (21)). Only the  $\{\gamma_t\}$  and some details associated with each  $w_t$  need to be stored and maintained. Let us now explain what these details are. It is easy to note that new  $w_t$  points arise only from  $\bar{z}$  in step 2. Thus each  $w_t$  has the representation,  $w_t = z_{i(t)} - z_{j(t)}$ , where  $i(t) \in I$  and  $j(t) \in J$ . The pair,  $(i(t), j(t))$  carries full information about  $w_t$  and hence it is sufficient to store this detail for each  $t$ .

Step 1 can be implemented using (22) and (23). To do the computations efficiently, it is useful to maintain a cache of the inner product,  $e_s = z \cdot z_s$ , for all those indices,  $s \in \{1, \dots, m\}$  which are needed in the definition of  $\{w_t\}$ . Let  $S$  denote the set of such indices, i.e.,  $S = \{i(t)\}_t \cup \{j(t)\}_t$ . In step 2,  $\tilde{z}$  takes the form,  $\tilde{z} = (1 - \lambda)z + \lambda\bar{z}$  where  $\bar{z} = z_i - z_j$  and  $i \in I, j \in J$  are determined in step 1. At the end of step 2, the inner products cache,  $\{e_s\}$  is updated using the following sequence of steps. (1) If  $i \notin S$ , set  $S := S \cup \{i\}$ , compute

$$e_i = z \cdot z_i = \sum_t \gamma_t w_t \cdot z_i = \sum_t \gamma_t (z_{i(t)} \cdot z_i - z_{j(t)} \cdot z_i)$$

and include it in the cache. (2) Repeat the above step for  $j$ . (3) Then for each  $s \in S$  set:  $e_s := (1 - \lambda)e_s + \lambda(z_i \cdot z_s - z_j \cdot z_s)$ . Then the  $\gamma_t$  are updated by the following sequence of steps. (1) Set:  $\gamma_t := (1 - \lambda)\gamma_t \forall t$ . (2) If  $z_i - z_j$  is equal to some already existing  $w_t$  then make a further

modification to  $\gamma_t$  (only for that  $t$ ):  $\gamma_t := \gamma_t + \lambda$ . (3) On the other hand, if  $z_i - z_j$  is not equal to any of the already existing  $w_t$  then include  $z_i - z_j$  as a new member of  $\{w_t\}$  and set its  $\gamma_t$  to  $\lambda$ . Note that the updatings involve two training vectors ( $z_i$  and  $z_j$ ) and require  $2n_t$  kernel evaluations where  $n_t$  is the number of  $t$  indices in the representation of  $z$ . Since  $n_t$  is large (hundreds or thousands) for most problems, the cost of computing the kernels and doing the updating of cache and  $\gamma_t$  is usually much higher than the cost of doing all other operations (such as minimizing norm on a line segment). A similar comment can be made for all algorithms considered in this paper.

At each iteration of the algorithm, a pair of points,  $u \in U$  and  $v \in V$  that yield  $z = u - v$  is always available:

$$u = \sum_t \gamma_t z_{i(t)} \quad v = \sum_t \gamma_t z_{j(t)}$$

Typically the condition,  $g_Z(-z, z) = 0$  does not occur at a finite  $k$ . However, since  $z \rightarrow z^*$  as the number of iterations goes to infinity, we have  $g_Z(-z, z) \rightarrow 0$ . It is easy to use (14) to see that

$$g_Z(-z, z) = g(u, v)$$

Because of this, the algorithm can be terminated (after a finite number of iterations) when

$$g_Z(-z, z) \leq \epsilon \|z\|^2$$

is satisfied, leading to the determination of an approximate solution of SVM-NV that satisfies the bounds given in Theorem 3.

We implemented Gilbert's algorithm using all the above ideas and tested its performance on a variety of classification problems. While the algorithm always makes rapid movement towards the solution during its initial iterations, on many problems it was very slow as it approached the final solution. Also, suppose there is a  $w_t$  which gets picked up by the algorithm during its initial stages, but which is not needed at all in representing the final solution (i.e.,  $z^* \cdot w_t > z^* \cdot z^*$ ), then the algorithm is slow in driving the corresponding  $\gamma_t$  to 0. Because of these reasons, Gilbert's algorithm, by itself, is not very efficient for solving the NPPs that arise in SVM classifier design.

## 5.2. Mitchell-Dem'yanov-Malozemov Algorithm.

Many years after Gilbert's work, Mitchell, Dem'yanov and Malozemov[18] independently suggested a new algorithm for MNP. We will refer to their algorithm as the MDM algorithm. Unlike Gilbert's algorithm, MDM algorithm fundamentally uses the representation,  $z = \sum_t \gamma_t w_t$  in its basic operation. When  $z = z^*$ , it is clear that  $z \cdot w_t = z \cdot z$  for all  $t$  which have  $\gamma_t > 0$ . If  $g_Z(-z, z) = 0$  then  $z = z^*$  and hence the above condition automatically holds. However, at an approximate point,  $z \neq z^*$ , the situation is different. The point,  $z$  could be very close to  $z^*$  and yet, there could very

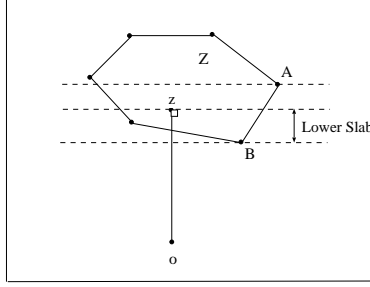


Figure 6: Definition of  $\Delta$ . Here:  $A = w_{tmin}$ ;  $B = s_Z(-z)$ ; and  $\Delta(z) = (-z \cdot B) - (-z \cdot A)$ .

well exist one or more  $w_t$  such that

$$z \cdot w_t \gg z \cdot z \quad \text{and} \quad \gamma_t > 0 \quad (24)$$

(though  $\gamma_t$  is small). For efficiency of representation as well as algorithm performance, it is a good idea to eliminate such  $t$  from the representation of  $z$ . With this in mind, define the function (see Figure 6),

$$\Delta(z) = h_Z(-z) - \min_{t:\gamma_t>0} (-z \cdot w_t) \quad (25)$$

Clearly,  $\Delta(z) \geq 0$ , and,  $\Delta(z) = 0$  if and only  $g_Z(-z, z) = 0$ , i.e.,  $z = z^*$ .

At each iteration, MDM algorithm attempts to decrease  $\Delta(z)$  whereas Gilbert's algorithm only tries to decrease  $g_Z(-z, z)$ . This can be seen geometrically in Figure 6. MDM algorithm tries to crush the total slab towards zero while Gilbert's algorithm only attempts to push the lower slab to zero. This is the essential difference between the two algorithms. Note that, if there is a  $t$  satisfying (24) then MDM algorithm has a much better potential to quickly eliminate such  $t$  from the representation.

Let us now describe the main step of MDM algorithm. Let  $tmin$  be an index such that

$$-z \cdot w_{tmin} = \min_{t:\gamma_t>0} (-z \cdot w_t)$$

Let  $d = s_Z(-z) - w_{tmin}$ . With the reduction of  $\Delta(z)$  in mind, MDM algorithm looks for improvement of norm value along the direction  $d$  at  $z$ . For  $\lambda > 0$ , let  $z = z + \lambda d$  and  $\phi(\lambda) = \|z\|^2$ . Since  $\phi'(0) = 2z \cdot d = -2\Delta(z)$ , it is easy to see that moving along  $d$  will strictly decrease  $\|z\|$  if  $z \neq z^*$ . Since, in the representation for  $z$ ,  $\gamma_{tmin}$  decreases during this movement,  $\lambda$  has to be limited to the interval defined by  $0 \leq \lambda \leq \gamma_{tmin}$ . Define  $\bar{z} = z + \gamma_{tmin}d$ . So, the basic iteration of MDM algorithm consists of finding the point of minimum norm on the line segment joining  $z$  and  $\bar{z}$ .

The structure of MDM algorithm is essentially the same as Gilbert's algorithm and so we feel there is no need to state it formally. Also, all the remarks that we made for Gilbert's algorithm

concerning its adaptation to the solution of NPPs arising in SVM classifier design hold good for MDM algorithm too. The details are only slightly more complicated because of the fact that the definition of  $\bar{z}$  involves 4 members of  $\{z_i\}$ . In particular, note that the updating of the  $\{e_s\}$  requires  $4n_t$  kernel evaluations in each iteration, where  $n_t$  is the number of  $t$  indices in the representation of  $z$ . We omit giving full details of the updates so as to avoid repetition.

We implemented and tested MDM algorithm. It works faster than Gilbert's algorithm, especially in the end stages when  $z$  approaches  $z^*$ . Algorithms which are much faster than MDM algorithm can be designed using the following two observations: (1) it is easy to combine the ideas of Gilbert's algorithm and MDM algorithm into a hybrid algorithm which is faster; and (2) by working directly in the space where  $U$  and  $V$  are located it is easy to find just two elements of  $\{z_k\}$  that are used to modify the  $z$ , and at the same time, keep the essence of the hybrid algorithm mentioned above. In the next section we describe the ideas behind the hybrid algorithm. Then we take up details of the final algorithm incorporating both observations, in section 6.

### 5.3. A Hybrid Algorithm.

A careful look at MDM algorithm shows that the main computational effort of each iteration is associated with the computation of  $h_Z(-z)$ , the determination of  $tmin$  and the updating of  $\gamma_t$  and of the inner products cache,  $e_s$ . Hence we can quite afford to make the remaining steps of an iteration a little more complex, provided that leads to some gain and does not disturb the determination of  $tmin$ ,  $tmax$  and the updatings. Since Gilbert's algorithm mainly requires the computation of  $h_Z(-z)$  and  $s_Z(-z)$ , which are anyway computed by the MDM algorithm, it is possible to easily insert the main idea of Gilbert's algorithm into MDM algorithm so as to improve it. A number of ideas emerge from this attempt. We will first describe these ideas and then comment on their goodness.

The situation at a typical iteration is shown in Figure 7. As before, let us take the representation for  $z$  as

$$z = \sum_t \gamma_t w_t$$

Without loss of generality let us assume that  $s_Z(-z) = w_{tmax}$  for some index,  $tmax$ . (If  $s_Z(-z)$  is not equal to any  $w_t$ , we can always include it as one more  $w_t$  and set the corresponding  $\gamma_t$  to 0 without affecting the representation of  $z$  in anyway.) The points  $w_{tmin}$  and  $w_{tmax}$  are respectively shown as  $A$  and  $B$  in the figure. The point,  $z + \gamma_{tmin}(w_{tmax} - w_{tmin})$  (i.e., the  $\bar{z}$  of MDM algorithm) is shown as  $C$ .

**Idea 1.** At step 2 of each iteration of MDM algorithm, take  $\tilde{z}$  to be the point of minimum norm on  $T_1$ , the triangle formed by  $z$ ,  $C$  and  $B$ .

An algorithm for computing the point of minimum norm on a triangle is given in the appendix.

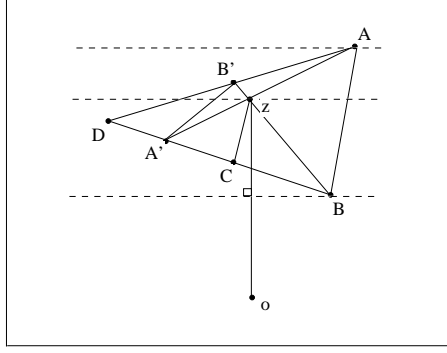


Figure 7: Illustrating the various ideas for improvement. Here:  $A = w_{tmin}$ ;  $B = s_Z(-z) = w_{tmax}$ ;  $zB$  = Gilbert line segment;  $zC$  = MDM line segment ( $zC$  is parallel to  $AB$ );  $T_1$  = triangle  $zCB$ ;  $T_2$  = triangle  $zA'B$ ;  $Q_3$  = quadrilateral  $ABA'B'$ ; and,  $T_4$  = triangle  $DAB$ .

Since the line segment connecting  $z$  and  $C$  (the one on which MDM algorithm finds the minimum point) is a part of  $T_1$ , this idea will locally perform (in terms of decreasing norm) at least as well as MDM algorithm. Note also that, since the minimum norm point on  $T_1$  can be expressed as a linear combination of  $z$ ,  $A$  and  $B$ , the cost of updating  $\gamma_t$  and the cache,  $e_s$  is the same as that of MDM algorithm.

For a given  $\bar{t}$ , let us define

$$z(\bar{t}) = \frac{1}{1 - \gamma_{\bar{t}}}(z - \gamma_{\bar{t}}w_{\bar{t}})$$

Clearly,  $z(\bar{t})$  is the point in  $Z$  obtained by removing  $w_{\bar{t}}$  from the representation of  $z$  and doing a re-normalization. Let us define:  $A' = z(tmin)$ ; see Figure 7. Since

$$C = z + \gamma_{tmin}(w_{tmax} - w_{tmin}) = (1 - \gamma_{tmin})A' + \gamma_{tmin}B$$

it follows that  $C$  lies on the line segment joining  $A'$  and  $w_{tmax}$ .

**Idea 2.** At step 2 of each iteration of MDM algorithm, take  $\tilde{z}$  to be the point of minimum norm on  $T_2$ , the triangle formed by  $z$ ,  $A'$  and  $B$ .

Because  $T_2$  contains  $T_1$ , idea 2 will produce a  $\tilde{z}$  whose norm is less than or equal to that produced by Idea 1.

One can carry this idea further to generate two more ideas. Let  $B' = z(tmax)$  and  $D$  be the point of intersection of the line joining  $A$ ,  $B'$  and the line joining  $B$ ,  $A'$ ; it is easily checked that  $D$  is the point obtained by removing both the indices,  $tmin$  and  $tmax$  from the representation of  $z$  and then doing a re-normalization, i.e.,

$$D = \frac{1}{(1 - \gamma_{tmin} - \gamma_{tmax})}(z - \gamma_{tmin}w_{tmin} - \gamma_{tmax}w_{tmax})$$

Let  $Q_3$  denote the quadrilateral formed by the convex hull of  $\{A, B, A', B'\}$  and  $T_4$  be the triangle formed by  $D, A$  and  $B$ . Clearly,  $T_2 \subset Q_3$ . Also, along the lines of an earlier argument, it is easy to show that  $Q_3 \subset T_4$ .

**Idea 3.** At step 2 of each iteration of MDM algorithm, take  $\tilde{z}$  to be the point of minimum norm on  $Q_3$ .

**Idea 4.** At step 2 of each iteration of MDM algorithm, take  $\tilde{z}$  to be the point of minimum norm on  $T_4$ .

We implemented and tested all of these ideas. Idea 1 gives a very good overall improvement over MDM algorithm and Idea 2 improves it further, a little more. However, we have found that, in overall performance, Idea 3 performs somewhat worse compared to Idea 2, and Idea 4 performs even worse! (We do not have a clear explanation for these performances.) On the basis of these empirical observations, we recommend Idea 2 to be the best one for use in modifying MDM algorithm.

## 6. A Fast Iterative Algorithm for NPP.

In this section we will give an algorithm for NPP directly in the space in which  $U$  and  $V$  are located. The key idea is motivated by considering Gilbert's algorithm. Let  $u \in U$  and  $v \in V$  be the approximate solutions at some given iteration and  $z = u - v$ .

We say that an index,  $k$  satisfies condition  $\mathcal{A}$  at  $(u, v)$  if:

$$k \in I \Rightarrow -z \cdot z_k + z \cdot u \geq \frac{\epsilon}{2} \|z\|^2 \quad (26)$$

$$k \in J \Rightarrow z \cdot z_k - z \cdot v \geq \frac{\epsilon}{2} \|z\|^2 \quad (27)$$

Suppose  $k$  satisfies (26). If  $\tilde{u}$  is the point on the line segment joining  $u$  and  $z_k$  that is closest to  $v$ , then, by the relation (A.8) given in the appendix we get an appreciable decrease in the objective function of NPP:

$$\|u - v\|^2 - \|\tilde{u} - v\|^2 \geq \min(\tilde{\epsilon}, \frac{\tilde{\epsilon}}{L^2}) \quad (28)$$

where  $\tilde{\epsilon} = 0.5\epsilon\|z\|^2$ . A parallel comment can be made if  $k$  satisfies (27). On the other hand, if we are unable to find an index  $k$  satisfying  $\mathcal{A}$  then (recall the definitions from section 3) we get:

$$g_U(-z, u) \leq \frac{\epsilon}{2} \|z\|^2, \quad g_V(z, v) \leq \frac{\epsilon}{2} \|z\|^2, \quad g(u, v) \leq \epsilon \|z\|^2 \quad (29)$$

which, by Theorem 3, is a good way of stopping. These observations lead us to the following generic algorithm, which gives ample scope for generating a number of specific algorithms from it. The efficient algorithm that we will describe soon after is one such special instance.

### Generic Iterative Algorithm for NPP.

0. Choose  $u \in U, v \in V$  and set  $z = u - v$ .

1. Find an index  $k$  satisfying  $\mathcal{A}$ . If such an index cannot be found, stop with the conclusion that the approximate optimality criterion, (29) is satisfied. Else go to step 2 with the  $k$  found.

2. Choose two convex polytopes,  $\tilde{U} \subset U$  and  $\tilde{V} \subset V$  such that  $u \in \tilde{U}$ ,  $v \in \tilde{V}$  and

$$z_k \in \tilde{U} \text{ if } k \in I, \quad z_k \in \tilde{V} \text{ if } k \in J \quad (30)$$

Compute  $(\tilde{u}, \tilde{v})$  to be a pair of closest points minimizing the distance between  $\tilde{U}$  and  $\tilde{V}$ . Set  $u = \tilde{u}$ ,  $v = \tilde{v}$  and go back to step 1.  $\square$

Suppose, in step 2 we do the following. If  $k \in I$ , choose  $\tilde{U}$  as the line segment joining  $u$  and  $z_k$ , and  $\tilde{V} = \{v\}$ . Else, if  $k \in J$ , choose  $\tilde{V}$  as the line segment joining  $v$  and  $z_k$ , and  $\tilde{U} = \{u\}$ . Then the algorithm is close in spirit to Gilbert's algorithm. Note however, that here only one index,  $k$  plays a role in the iteration, whereas Gilbert's algorithm requires two indices (one each from  $I$  and  $J$ ) to define  $\bar{z}$ . In a similar spirit, by appropriately choosing  $\tilde{U}$  and  $\tilde{V}$ , the various ideas of section 5.3 can be extended while involving only two indices. Before we discuss a detailed algorithm, we prove convergence of the generic algorithm.

**Theorem 5.** The generic iterative algorithm terminates in step 1 with a  $(u, v)$  satisfying (29), after a finite number of iterations.  $\square$

The proof is easy. First, note that (28) holds because of the way  $\tilde{U}$  and  $\tilde{V}$  are chosen in step 2. Since, by assumption A1,  $\|z\|^2$  is uniformly bounded below by zero, there is a uniform decrease in  $\|z\|^2$  at each iteration. Since the minimum distance between  $U$  and  $V$  is nonnegative the iterations have to terminate within a finite number of iterations.  $\square$

Consider the situation in step 2 where an index  $k$  satisfying  $\mathcal{A}$  is available. Suppose the algorithm maintains the following representations for  $u$  and  $v$ :

$$u = \sum_{i \in \hat{I}} \beta_i z_i, \quad v = \sum_{j \in \hat{J}} \beta_j z_j \quad (31)$$

where  $\hat{I} \subset I$ ,  $\hat{J} \subset J$ , and

$$\beta_k > 0 \quad \forall k \in \hat{I} \cup \hat{J}, \quad \sum_{i \in \hat{I}} \beta_i = 1, \quad \sum_{j \in \hat{J}} \beta_j = 1$$

We will refer to a  $z_k$ ,  $k \in \hat{I} \cup \hat{J}$  as a *support vector*, consistent with the terminology adopted in the literature. Let us define  $\text{imin}$  and  $\text{jmin}$  as follows (see Figure 8):

$$\begin{aligned} \text{imin} &= \arg \min \{-z \cdot z_i : i \in \hat{I}, \beta_i > 0\} \\ \text{jmin} &= \arg \min \{z \cdot z_j : j \in \hat{J}, \beta_j > 0\} \end{aligned} \quad (32)$$

(Though redundant, we have stated the conditions,  $\beta_i > 0$  and  $\beta_j > 0$  to stress their importance. Also, in a numerical implementation, an index  $k \in \hat{I} \cup \hat{J}$  having  $\beta_k = 0$  could occur numerically.

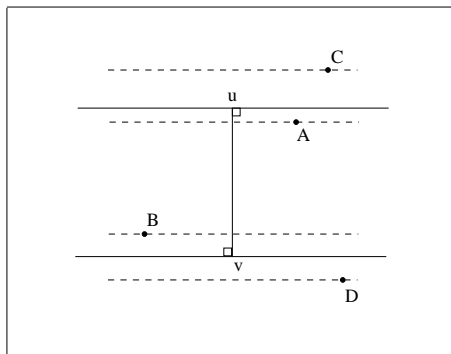


Figure 8: Illustration of  $i_{\min}$ ,  $j_{\min}$ ,  $i_{\max}$  and  $j_{\max}$ . Here  $C = z_{i_{\min}}$ ,  $D = z_{j_{\min}}$ ,  $A = z_{i_{\max}}$  and  $B = z_{j_{\max}}$ .

Unless such indices are regularly cleaned out of  $\hat{I} \cup \hat{J}$ , it is best to keep these positivity checks in (32.)

There are a number of ways of combining  $u$ ,  $v$ ,  $z_k$  and one index from  $\{i_{\min}, j_{\min}\}$  and then doing operations along ideas similar to those in section 5.3 to create a variety of possibilities for  $\tilde{U}$  and  $\tilde{V}$ . We have implemented and tested some of the promising ones and empirically arrived at one final choice, which is the only one that we will describe in detail.

First we set  $k_{\min}$  as follows:  $k_{\min} = i_{\min}$  if  $-z \cdot z_{i_{\min}} + z \cdot u < z \cdot z_{j_{\min}} - z \cdot v$ ; else  $k_{\min} = j_{\min}$ . Then step 2 is carried out using one of the following four cases, depending on  $k$  and  $k_{\min}$ . See Figure 9 for a geometric description of the cases.

*Case 1.*  $k \in I$ ,  $k_{\min} \in I$ . Let  $C'$  be the point in  $U$  obtained by removing  $z_{k_{\min}}$  from the representation of  $u$ , i.e.,

$$C' = \frac{1}{(1 - \beta_{k_{\min}})}(u - \beta_{k_{\min}}z_{k_{\min}}) = u + \mu(u - z_{k_{\min}}) \quad (33)$$

where  $\mu = \beta_{k_{\min}}/(1 - \beta_{k_{\min}})$ . Choose:  $\tilde{U}$  to be the triangle formed by  $u$ ,  $C'$  and  $z_k$ ; and,  $\tilde{V} = \{v\}$ .

*Case 2.*  $k \in J$ ,  $k_{\min} \in J$ . Let  $D'$  be the point in  $V$  obtained by removing  $z_{k_{\min}}$  from the representation of  $v$ , i.e.,

$$D' = \frac{1}{(1 - \beta_{k_{\min}})}(v - \beta_{k_{\min}}z_{k_{\min}}) = v + \mu(v - z_{k_{\min}}) \quad (34)$$

where  $\mu = \beta_{k_{\min}}/(1 - \beta_{k_{\min}})$ . Choose:  $\tilde{V}$  to be the triangle formed by  $v$ ,  $D'$  and  $z_k$ ; and,  $\tilde{U} = \{u\}$ .

*Case 3.*  $k \in I$ ,  $k_{\min} \in J$ . Choose:  $\tilde{U}$  to be the line segment joining  $u$  and  $z_k$ ; and  $\tilde{V}$  to be the line segment joining  $v$  and  $D'$  where  $D'$  is as in (34).

*Case 4.*  $k \in J$ ,  $k_{\min} \in I$ . Choose:  $\tilde{V}$  to be the line segment joining  $v$  and  $z_k$ ; and  $\tilde{U}$  to be the line segment joining  $u$  and  $C'$  where  $C'$  is as in (33).

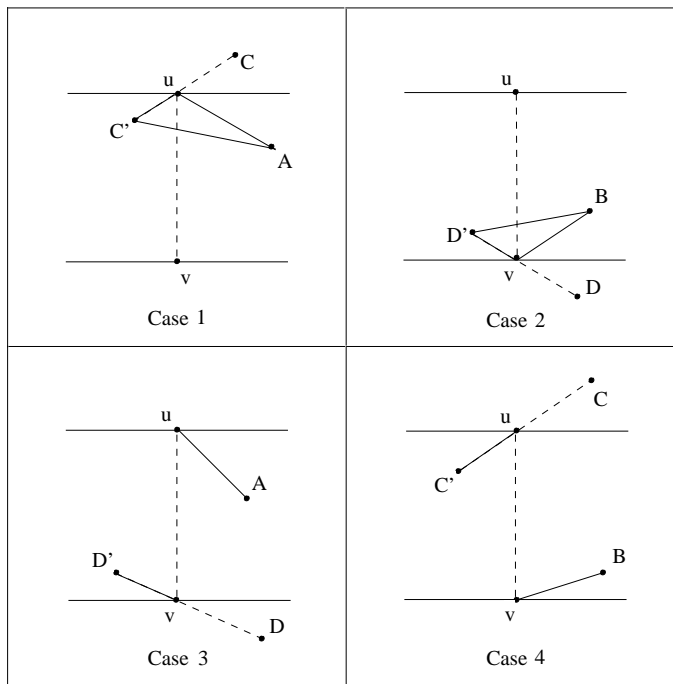


Figure 9: A geometric description of the 4 cases. For Cases 1 and 3,  $A = z_k$ ; for Cases 2 and 4,  $B = z_k$ .

This defines the basic operation of the algorithm. Efficient algorithms for doing nearest point computations involving a triangle and line segments can be found in the appendix and [17]. Tremendous improvement in efficiency can be achieved by doing two things: (1) maintaining and updating caches for some variables; and (2) interchanging operations between the support vector and non-support vector sets. These ideas are directly inspired from those used by Platt in his SMO algorithm. Let us now go into these details.

First let us discuss the need for maintaining cache for some variables. A look at (26), (27) and (32) shows that  $u \cdot z_k$ ,  $v \cdot z_k$ ,  $u \cdot u$ ,  $u \cdot v$ ,  $v \cdot v$  and  $z \cdot z$  are important variables. If any of these variables is to be computed from scratch, it is expensive; for example, computation of  $u \cdot z_k$  (for a single  $k$ ) using (31) requires the evaluation of  $m_1$  kernel computations and  $O(m_1)$  computation (where  $m_1 = |\hat{I}|$ , the number of elements in  $\hat{I}$ ) and, to compute  $\text{imin}$  of (32) this has to be repeated  $m_1$  times! Therefore it is important to maintain caches for these variables. Let us define:  $e_u(k) = u \cdot z_k$ ,  $e_v(k) = v \cdot z_k$ ,  $\delta_{uu} = u \cdot u$ ,  $\delta_{vv} = v \cdot v$ ,  $\delta_{uv} = u \cdot v$  and  $\delta_{zz} = z \cdot z$ .

Among these cache variables,  $e_u$  and  $e_v$  are the only vectors. Since the algorithm spends much of its operation adjusting the  $\beta$  of the support vectors, it is appropriate to bring in the non-support vectors only rarely to check satisfaction of optimality criteria. In that case, it is efficient to maintain  $e_u(k)$  and  $e_v$  only for  $k \in \hat{I} \cup \hat{J}$  and, for other indices compute these quantities from scratch whenever

needed.

As mentioned in the last paragraph, it is important to spend most of the iterations using  $k \in \hat{I} \cup \hat{J}$ . With this in mind, we define two types of loops. The *Type I* loop goes *once* over all  $k \notin \hat{I} \cup \hat{J}$ , sequentially, choosing one at a time, checking condition  $\mathcal{A}$ , and, if it is satisfied, doing one iteration, i.e., step 2 of the generic algorithm. The *Type II* loop operates only with  $k \in \hat{I} \cup \hat{J}$ , doing the iterations many many times until certain criteria are met. Whereas, in Platt's SMO algorithm, the iterations run over one support vector index at a time, we have adopted a different strategy for Type II loop. Let us define  $\text{imax}$ ,  $\text{jmax}$  and  $\text{kmax}$  as follows.

$$\begin{aligned} \text{imax} &= \arg \max\{-z \cdot z_i : i \in \hat{I}\} \\ \text{jmax} &= \arg \max\{z \cdot z_j : j \in \hat{J}\} \end{aligned} \quad (35)$$

$$\text{kmax} = \begin{cases} \text{imax} & \text{if } -z \cdot z_{\text{imax}} + z \cdot u > z \cdot z_{\text{jmax}} - z \cdot v \\ \text{jmax} & \text{otherwise} \end{cases} \quad (36)$$

These can be efficiently computed mainly because of the caching of  $e_u(k)$  and  $e_v(k)$  for all  $k \in \hat{I} \cup \hat{J}$ . A basic Type II iteration consists of determining  $\text{kmax}$  and doing one iteration using  $k = \text{kmax}$ . These iterations are continued until the approximate optimality criteria are satisfied over the support vector set, i.e., when

$$\begin{aligned} g_U^{\text{supp}} &= -z \cdot z_{\text{imax}} + z \cdot u \leq \frac{\epsilon}{2} \|z\|^2 \\ g_V^{\text{supp}} &= z \cdot z_{\text{jmax}} - z \cdot v \leq \frac{\epsilon}{2} \|z\|^2 \end{aligned} \quad (37)$$

In the early stages of the algorithm (say a few of the TypeI-then-TypeII rounds) a “nearly accurate” set of support vector indices get identified. Until that happens, we have found that it is wasteful to spend too much time doing Type II iterations. Therefore we have adopted the following strategy. If, at the point of entry to a Type II loop, the percentage difference of the number of support vectors as compared with the value at the previous entry is greater than 2%, then we limit the number of Type II iterations to  $m$ , the total number of training pairs in the problem. This is done to roughly make the cost of Type I and Type II loops equal. In any case, we have found it useful to limit the number of Type II iterations to  $10m$ . The algorithm is stopped when (37) holds, causing the algorithm to exit Type II loop, and, in the Type I loop that follows there is no  $k$  satisfying condition  $\mathcal{A}$ .

The main ideas relating to Type II iterations that we have described above can also be used in Platt's SMO algorithm. We have done this and have also done some initial testing. From this limited testing, it appears that these ideas do not cause a significant difference to the SMO algorithm. On the other hand, these ideas make a big difference for our algorithm, i.e., implementing Type II by sequentially running  $k$  over all support vector indices many times progresses much much slower compared to the implementation based on  $\text{kmax}$ .

Full details concerning the actual implementation of the full algorithm can be found in [14].

## 7. A Simplification Using the $b^2$ term.

In this section we consider a simpler SVM formulation suggested recently by Friess et. al.[8] and Friess[9]. Suppose the term,  $b^2/2$  is added to the cost functions of SVM-NV and SVM-VQ that were formulated in section 1. Then by appropriate rearrangement these problems can be written as a special SVM-NV problem in which  $b$  is absent. Let us take SVM-VQ to illustrate this point. Similar to (1) we can define

$$\tilde{w} = \begin{pmatrix} w \\ \sqrt{C}\xi \\ b \end{pmatrix}; \quad \tilde{z}_i = \begin{pmatrix} z_i \\ \frac{1}{\sqrt{C}}e_i \\ 1 \end{pmatrix}, i \in I; \quad \tilde{z}_j = \begin{pmatrix} z_j \\ -\frac{1}{\sqrt{C}}e_j \\ 1 \end{pmatrix}, j \in J. \quad (35)$$

Then it is easy to see that SVM-V transforms to *an instance* of the following problem:

$$\begin{aligned} \min \quad & \frac{1}{2}\|w\|^2 \\ \text{s.t.} \quad & w \cdot z_i \geq 1 \quad \forall i \in I; \quad w \cdot z_j \leq -1 \quad \forall j \in J \end{aligned} \quad (\text{SVM - NV - Nob})$$

As in the discussion below equation (1), if  $K$  denotes the Kernel function in the original problem, then  $\tilde{K}$ , the Kernel function for the transformed problem is given by

$$\tilde{K}(x_k, x_l) = K(x_k, x_l) + \frac{1}{C}\delta_{kl} + 1$$

In the case of SVM-NV, when  $b^2/2$  is added to the objective function, the transformation to an instance of SVM-NV-Nob is easier than the above steps; we only have to leave out  $\xi$ .

The rest of this section will correspond to the notations in SVM-NV-Nob. Let us assume the following.

**Assumption A2.** The feasible space of SVM-NV-Nob is non-empty.  $\square$

As noted in Remark 2, SVM-VQ always has non-empty feasible space and hence, when it is transformed to SVM-NV-Nob, it also has a non-empty feasible space and the above assumption holds. Let  $W = \{z_i : i \in I\} \cup \{-z_i : i \in J\}$  and  $Z = \text{co}W$ . It is also easy to show that this problem is equivalent to:

$$\min \|z\| \quad \text{s.t.} \quad z \in Z \quad (MNP)$$

Parallel to Theorem 1 we can derive the following result.

**Theorem 6.**  $w^*$  solves SVM-NV-Nob if and only if  $z^*$  solves MNP and

$$w^* = \frac{1}{\|z^*\|^2} z^* \quad (38)$$

As in section 2, the following Wolfe dual problem of SVM-NV-Nob can be used to understand this result:

$$\max \sum_k \alpha_k - \frac{1}{2} \sum_k \sum_l \alpha_k \alpha_l y_k y_l z_k \cdot z_l \quad \text{s.t.} \quad \alpha_k \geq 0 \quad \forall k \quad (39)$$

Introduce  $\lambda = \sum_k \alpha_k$ ,

$$\beta_k = \frac{\alpha_k}{\lambda} \quad \forall k \quad (40)$$

to rewrite (39) as

$$\begin{aligned} \max \quad & \lambda - \frac{\lambda^2}{2} \sum_k \sum_l \beta_k \beta_l y_k y_l z_k \cdot z_l \\ \text{s.t.} \quad & \beta_k \geq 0 \quad \forall k; \quad \sum_k \beta_k = 1. \end{aligned} \quad (41)$$

Optimizing  $\lambda$  for fixed  $\beta$  gives

$$\lambda^* = \frac{1}{\sum_k \sum_l \beta_k \beta_l y_k y_l z_k \cdot z_l} \quad (42)$$

Substituting this in (41) yields the equivalent problem

$$\begin{aligned} \min \quad & \frac{1}{2} \sum_k \sum_l \beta_k \beta_l y_k y_l z_k \cdot z_l \\ \text{s.t.} \quad & \beta_k \geq 0 \quad \forall k; \quad \sum_k \beta_k = 1. \end{aligned}$$

which is equivalent to MNP if we observe that  $P\beta \in \text{co}\{y_1 z_1, \dots, y_m z_m\}$  where  $P$  is as in the discussion below Theorem 1. Since  $w^* = \sum_k \alpha_k^* y_k z_k$ , where  $\alpha^*$  denotes a solution of (39), we can use (40) and (42) to easily derive (38).

**Remark 4.** Unlike (21), here the set  $W$  that defines  $Z$  is directly prescribed in the space where  $U$  and  $V$  are located, and that too, in a very simple fashion. Therefore, Gilbert's algorithm, MDM algorithm, and its modifications apply directly without need for Minkowski set difference transformations. Thus their implementation becomes much simpler than that described in section 6. Since the ideas are very straight forward, we do not go into the details.  $\square$

## 8. Computational Comparisons.

In this section we compare the performance of our algorithm(s) against other iterative algorithms for SVM classifier design. Since the objective of this paper is to design fast algorithms, our focus in this section is only on studying computational efficiency. Whether or not the various SVM formulations tried in this paper lead to good generalizing solutions is a different issue that is not dealt with here; it is important, however, to do a separate study of this issue.

The following 6 methods were chosen for comparison.

1. **SMO** : Platt's SMO was applied to SVM-VL.
2. **SMO\_Q** : SVM-VQ was converted to SVM-NV, which was then solved using SMO.

3. **NPA** : SVM-VQ was converted to SVM-NV, which was then solved by the near point algorithm of section 6.
4. **SOR** : Mangasarian and Musicant’s SOR was applied to SVM-VL-BSQ, which is the modification of SVM-VL obtained by adding the cost,  $b^2/2$  in the objective function.
5. **MNA** : SVM-VQ-BSQ, the modification of SVM-VQ obtained by adding the cost,  $b^2/2$  in the objective function, was converted to SVM-NV-Nob and then solved using the minimum norm algorithm that we briefly outlined at the end of section 7.
6. **SOR\_Q** : The same problem solved by MNA is solved using SOR. As we already mentioned in section 1, this way of using SOR is identical to Friess’ Adatron algorithm[9]. However we have preferred to call it as SOR\_Q since the heuristics (using two types of loops, sorting support vectors etc.) that cause significant improvements are basically taken from Mangasarian and Musicant’s paper[19].

We implemented all these methods in Fortran and ran them using *g77* on a 200 MHz Pentium machine. In the case of SMO and SOR, they were implemented exactly along the lines suggested by Platt[20] and Mangasarian and Musicant[19].

When testing an algorithm, it is important to specify clearly what criteria are used to stop the algorithm. When comparing different algorithms, this is even more important; also, in this case, it is crucial to select stopping criteria that are as uniform as possible. In our testing we have used stopping criteria which were suggested by Platt[20] and Joachims[12]. In each of the *transformed formulations* used by the six methods listed above, we have a pair of parallel hyperplanes that form the “boundaries” for the two classes. Let us call these hyperplanes as  $B_1$  and  $B_2$ . The definition of  $B_1$  and  $B_2$  differs for different transformed formulations. Let us define:

$$w = \sum_k \alpha_k y_k z_k, \quad \bar{u} = \sum_{i \in I} \beta_i z_i, \quad \bar{v} = \sum_{j \in J} \beta_j z_j, \quad \bar{z} = \sum_k \beta_k y_k z_k$$

Using these, Table 1 defines  $B_1$  and  $B_2$  for the six methods. Note that the notations in the last two columns correspond to the transformed formulation. For instance, for NPA, the  $z_k$ ’s used in the definitions correspond to those in SVM-NV, which are actually the  $\tilde{z}_k$ ’s defined in (1).

Let  $D$  denote the distance between the two hyperplanes,  $B_1$  and  $B_2$ . Given a relative tolerance,  $\epsilon$ , let us form symmetric bands, each of thickness,  $\epsilon D$ , using the four shifted hyperplanes,  $H_1$ ,  $H_2$ ,  $H_3$  and  $H_4$ , as illustrated in Figure 10. Let us simply refer to a variable used in the problems ( $\alpha_k$  or  $\beta_k$ ) as a *multiplier*. As we know, all multipliers have to satisfy non-negativity constraints. Further, in the case of SMO and SOR, the multipliers also have to satisfy the upper bound constraint:  $\alpha_k \leq C$ ; for notational purposes we will take  $C = \infty$  for all other methods. At a particular

| Method | Problem addressed | Transformed Formulation | Iteration Variables | Eqns of $B_1$ and $B_2$ in $z$ -space  | $D$            |
|--------|-------------------|-------------------------|---------------------|--|----------------|
| SMO    | SVM-VL            | No Transformation       | $\alpha, b$         | $w \cdot z + b = \pm 1$  | $2/\ w\ $      |
| SMO_Q  | SVM-VQ            | SVM-NV                  | $\alpha, b$         | $w \cdot z + b = \pm 1$  | $2/\ w\ $      |
| NPA    | SVM-VQ            | SVM-NV                  | $\beta$             | $\bar{z} \cdot z = \bar{z} \cdot \bar{u}, \bar{z} \cdot z = \bar{z} \cdot \bar{v}$ | $\ \bar{z}\ $  |
| SOR    | SVM-VL-BSQ        | SVM-VL without $b$      | $\alpha$            | $w \cdot z = \pm 1$  | $2/\ w\ $      |
| MNA    | SVM-VQ-BSQ        | SVM-NV-Nob              | $\beta$             | $\bar{z} \cdot z = \pm \bar{z} \cdot \bar{z}$                                      | $2\ \bar{z}\ $ |
| SOR_Q  | SVM-VQ-BSQ        | SVM-NV-Nob              | $\alpha$            | $w \cdot z = \pm 1$  | $2/\ w\ $      |

Table 1: Definitions of  $B_1$ ,  $B_2$  and  $D$ .

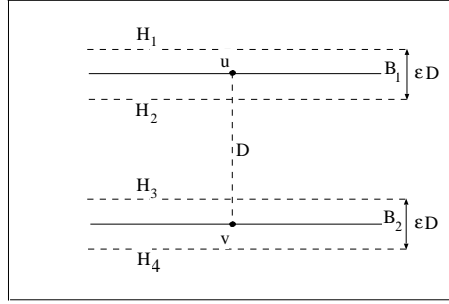


Figure 10: Illustration of the notations used in the stopping criteria.

instance we say that the  $k$ -th multiplier *violates* the stopping criteria if *any one* of the following four conditions is met.

- $y_k = 1$ , multiplier is positive and the resulting SVM output is above  $H_1$ .
- $y_k = -1$ , multiplier is positive and the resulting SVM output is below  $H_4$ .
- $y_k = 1$ , multiplier is lower than  $C$  and the resulting output is below  $H_2$ .
- $y_k = -1$ , multiplier is lower than  $C$  and the resulting output is above  $H_3$ .

An algorithm is stopped when there is no multiplier that violates the stopping criteria. *The value,  $\epsilon = 0.001$  was used for all experiments.* In the case of NPA, note that, for the same  $\epsilon$ , the stopping criteria just described is tighter than those suggested in section 4 based on the  $g$  function, i.e., (15). For the sake of uniform comparison, we have used the tighter criteria mentioned above.

The following standard problems were used in our testing: Wisconsin Breast Cancer data[26,4]; Two Spirals data[24]; Checkers data[13]; UCI Adult data[21]; and Web page classification data[21, 12]. Except for Checkers data, for which we created a random set of points on a  $4 \times 4$  checkers grid, all other data sets were downloaded from the internet sites mentioned in the above references and

| Data Set                | $\sigma^2$ | $n$ | $m$   |
|-------------------------|------------|-----|-------|
| Wisconsin Breast Cancer | 4.0        | 9   | 683   |
| Two Spirals             | 0.5        | 2   | 195   |
| Checkers                | 0.5        | 2   | 465   |
| Adult-1                 | 10.0       | 123 | 1605  |
| Adult-4                 | 10.0       | 123 | 4781  |
| Adult-7                 | 10.0       | 123 | 16100 |
| Web-1                   | 10.0       | 300 | 2477  |
| Web-4                   | 10.0       | 300 | 7366  |
| Web-7                   | 10.0       | 300 | 24692 |

Table 2: Data Set Properties.

were used in full for training, i.e., no division into training/validation/test sets was made. In the case of Adult data set, the inputs are represented in a special binary format, as used by Platt in the testing of SMO. To study scaling properties as training data grows, Platt did staged experiments on the Adult and Web data. We have used only the data from the first, fourth and seventh stages. The gaussian kernel,

$$K(x_i, x_j) = \exp(-0.5\|z_i - z_j\|^2/\sigma^2)$$

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 2. The  $\sigma^2$  values given in this table were chosen as follows. For the Adult and Web data the  $\sigma^2$  values are the same as those used by Platt in his experiments on SMO; for other data, we chose  $\sigma^2$  suitably to get good generalization.

Mangasarian and Musicant pointed out that the final solutions obtained with and without the  $b^2/2$  term included in the cost function of an SVM formulation are typically close. We also have observed this in our experiments. Given this, the methods, SMO\_Q, NPA, MNA and SOR\_Q can be directly compared; similarly, SMO and SOR can be directly compared. To make a cross comparison between these two groups of methods, we did the following. Let  $M$  denote the final margin obtained by a solution in the solution space. For each data set, we chose two different ranges for  $C$  and  $\tilde{C}$  values in such a way that they cover the same range of  $M$  values. (It may be noted that  $M$  has an inverse relationship with  $C$  and  $\tilde{C}$ .) Then we ran all methods on a bunch of  $C$  and  $\tilde{C}$  values sampled from those ranges. Such a study is important for another reason, as well. When a particular method is used for SVM design,  $C$  or  $\tilde{C}$  is usually unknown, and it has to be chosen by trying a number of values and using a validation set. Therefore, good performance of a method in a range of  $C$  or  $\tilde{C}$  values is important. Whether or not a particular method has such a performance gets revealed clearly in our experiments.

We use the total number of kernel evaluations as the measure of the total computing cost. Since the total number of kernel evaluations is pretty much independent of the computing environment used, it is easy for other researchers developing new algorithms to compare their methods against the ones investigated here, without actually running these methods again. Our argument for using total number of kernel evaluations to measure computational cost is as follows. For each of the methods, the cost of doing one successful step, `Cost_Per_Step` can be written as:

$$\text{Cost\_Per\_Step} = C_o + n_{sv} * [C_i + n_k * C_k + C_e]$$

where:  $C_o$  = overhead cost in each step;  $n_{sv}$  = number of support vectors;  $n_k$  = number of kernel evaluations per support vector in each step;  $C_k$  = cost of doing one kernel evaluation; and the following are costs per support vector index:  $C_i$  = cost of choosing indices to use in the step;  $C_e$  = cost of updating caches after the kernel evaluations are done. Rewriting the expression for `Cost_Per_Step` we get:

$$\text{Cost\_Per\_Step} = n_{sv} * n_k * [C_o / (n_{sv} * n_k) + C_i / n_k + C_k + C_e / n_k]$$

To illustrate our point, let us only compare SOR, SMO and NPA. Measuring cost by number of arithmetic operations,  $C_o$  is around 7 for SOR, 30 for SMO and 120 for NPA.  $C_i$  is 0 for SOR, 2 for SMO (needed for choosing second choice heuristic) and 4 for NPA (needed for choosing two indices).  $C_e/n_k$  is same for all three methods. For SOR  $n_k$  is 1, while  $n_k$  is 2 for both, SMO and NPA. In the case of dense kernel evaluation,  $C_k$  is easily more than  $2n$  operations, where  $n$  is the input dimension. Even if we take  $n$  to be in the order of 5 and  $n_{sv}$  to be in the order of 100, it is clear that  $C_o / (n_{sv} * n_k)$  and  $C_i / n_k$  are negligible when compared to  $C_k$ . Even in the case of sparse kernel computation, the cost,  $C_k$  is typically dominant (since a number of integer comparisons are involved). Because of these reasons, we feel  $n_{sv} * n_k$  is a good indicator of `Cost_Per_Step` and hence, the total number of kernel evaluations needed for solving a problem is a good indicator of the total computing time. We note, however, that in the case of linear SVMs, the situation is quite different and `Cpu_Time` has to be actually measured to do a fair comparison of the methods. We are also assuming that typical problems are large enough not to allow the storage of the entire kernel matrix and hence the kernel values are to be computed as and when needed. We haven't studied the possibility of using a cache for kernels; when such a cache is possible, it is expected that the efficiency of each method will be enhanced differently and comparison has to be done using `Cpu_Times`.

For the various test problems, figures 11-21 show the variation of computing cost as a function of the solution space margin,  $M$ . `SOR_Q` does not appear in many figures. This is because computing times involved in these runs were very large and so `SOR_Q` was not run till completion. In the case of Adult-7 and Web-7 data we compared only SMO and NPA.

The following conclusions emerge from our computational experiments.

1. For the SVM problem formulation which linearly penalizes classification violations, SMO is clearly better than SOR. Although the experiments reported in [19] point to SMO and SOR to be equally efficient on large datasets, those experiments were for Linear SVMs with  $C$  fixed to a single value. The comparisons reported in this paper are for Nonlinear SVMs using Gaussian kernels. If one looks at the variation of computing cost as a function of  $C$ , (see, for example, the performances on adult and web datasets) it can be noted that SOR's cost becomes high as  $C$  increases. Even for low  $C$  values SMO is quite better.
2. For the SVM problem formulation which quadratically penalizes classification violations, NPA is clearly the best.
3. SOR\_Q (Adatron) is a very slow method. However it must be noted that when the storage of the entire kernel matrix is possible, methods such as SOR\_Q (Adatron) can be competitive. Also, the experiments in [8,9] point out that good generalizing solutions are usually obtained within a small number of loops of the training set. Hence these methods are still valuable for upto medium-sized problems.
4. SOR and MNA, the other two methods based on inclusion of  $b^2/2$  term in the cost function, perform reasonably well, but not as fast as SMO and NPA. Also, their computational cost increases as  $C, \tilde{C}$  are increased to large values.
5. When SMO and NPA are compared over a number of  $C, \tilde{C}$  values, they perform quite closely, on the average. The runs on Adult and Web data point to the fact that both these methods scale equally well as a function of data size. At large  $C, \tilde{C}$  values, NPA seems to have an edge over SMO on most datasets, which indicates that it is the better method for solving SVM-NV. Overall these two methods give the best performance.
6. The linear violation formulation results in a smaller number of support vectors when compared to the quadratic violation formulation. While the difference is very little in the Wisconsin Breast Cancer, Two Spirals, and Checkers datasets, the difference is very prominent in Adult and Web datasets, especially at small  $C, \tilde{C}$  values. If such a large difference does occur in a particular problem, then the linear formulation does have a great advantage when using the SVM classifier for classification after the design process is over.

## 9. Conclusion.

In this paper we have given a new, fast algorithm for SVM classifier design using a carefully devised nearest point algorithm. The performance of this algorithm on a number of benchmark problems is very encouraging. Though not as simple as the SMO algorithm, our algorithm also

is quite straightforward to implement, as indicated by the pseudocode in [14]. Our nearest point formulation as well as the algorithm are special to classification problems and cannot be used for SVM regression.

## References

- [1] J.K. Anlauf and M. Biehl, The Adatron: an adaptive perceptron algorithm, *Europhysics Letters*, Vol.10, No.7, 1989, pp.687-692.
- [2] R.O. Barr, An efficient computational procedure for a generalized quadratic programming problem, *SIAM J. Control*, Vol.7, No.3, 1969, pp.415-429.
- [3] R. Bennett and E.J. Bredensteiner, Geometry in learning, Tech. Report, Department of Mathematical Sciences, Rennselaer Polytechnic Institute, New York, 1996.
- [4] R. Bennett and O.L. Mangasarian, Robust linear programming discrimination of two linearly inseparable sets, *Optimization Methods and Software*, Vol.1, 1992, pp.23-34.
- [5] C.J.C. Burges, A tutorial on support vector machines for pattern recognition, *Data Mining and Knowledge Discovery*, Vol.2, Number 2, 1998.
- [6] C. Cortes and V. Vapnik, Support Vector Networks, *Machine Learning*, Vol.20, 1995, pp.273-297.
- [7] R. Fletcher, *Practical Methods of Optimizations*. John Wiley and Sons, 2nd Edition, 1987.
- [8] T.T. Friess, N. Cristianini, and C. Campbell, The kernel adatron algorithm: a fast and simple learning procedure for support vector machines, *Proc. 15th International Conference on Machine Learning*, Morgan Kaufman Publishers, 1998.
- [9] T.T. Friess, Support vector networks: The kernel adatron with bias and soft-margin, Tech. Report, The University of Sheffield, Dept. of Automatic Control and Systems Engineering, Sheffield, England, 1998.
- [10] E.G. Gilbert, Minimizing the quadratic form on a convex set, *SIAM J. Control*, Vol.4, 1966, pp.61-79.
- [11] E.G. Gilbert, D.W. Johnson and S.S. Keerthi, A fast procedure for computing the distance between complex objects in three dimensional space, *IEEE J. Robotics and Automation*, Vol.4, 1988, pp.193-203.

- [12] T. Joachims, Making large-scale support vector machine learning practical, in B. Schölkopf, C. Burges, A. Smola. *Advances in Kernel Methods: Support Vector Machines*, MIT Press, Cambridge, MA, December 1998.
- [13] L. Kaufman, Solving the quadratic programming problem arising in support vector classification, in B. Schölkopf, C. Burges, A. Smola. *Advances in Kernel Methods: Support Vector Machines*, MIT Press, Cambridge, MA, December 1998.
- [14] S.S. Keerthi, S.K. Shevade, C. Bhattacharyya and K.R.K. Murthy, A fast iterative nearest point algorithm for support vector machine classifier design, Tech. Report TR-ISL-99-03, Intelligent Systems Lab, Dept. of Computer Science and Automation, Indian Institute of Science, Bangalore, India, March 1999. See: <http://guppy.mpe.nus.edu.sg/~mpessk>
- [15] C.L. Lawson and R.J. Hanson, *Solving least squares problems*. Prentice Hall, 1974.
- [16] S.R. Lay, *Convex sets and their applications*. John Wiley and Sons, 1982.
- [17] V.J. Lumelsky, On fast computation of distance between line segments, *Information Processing Letters*, Vol.21, 1985, pp.55-61.
- [18] B.F. Mitchell, V.F. Dem'yanov, and V.N. Malozemov, Finding the point of a polyhedron closest to the origin, *SIAM J. Control*, Vol.12, 1974, pp.19-26.
- [19] O.L. Mangasarian and D.R. Musicant, Successive overrelaxation for support vector machines, Tech. Report, Computer Sciences Dept., University of Wisconsin, Madison, WI, USA, 1998.
- [20] J.C. Platt, Fast training of support vector machines using sequential minimal optimization, in B. Schölkopf, C. Burges, A. Smola. *Advances in Kernel Methods: Support Vector Machines*, MIT Press, Cambridge, MA, December 1998.
- [21] J.C. Platt, Adult and Web Datasets. <http://www.research.microsoft.com/~jplatt>
- [22] N.K. Sancheti and S.S. Keerthi, Computation of certain measures of proximity between convex polytopes: A complexity viewpoint, *Proc. of IEEE International Conference on Robotics and Automation*, Nice, France, 1992, pp.2508-2513.
- [23] A.J. Smola and B. Schölkopf, A tutorial on support vector regression, *NeuroCOLT Technical Report TR-1998-030*, Royal Holloway College, London, UK, 1998.
- [24] Two Spirals Data.  
<ftp://ftp.boltz.cs.cmu.edu/pub/neural-bench/bench/two-spirals-v1.0.tar.gz>
- [25] V. Vapnik, *The Nature of Statistical Learning Theory*. Springer-Verlag, New York, 1995.

[26] Wisconsin Breast Cancer Data.

<ftp://128.195.1.46/pub/machine-learning-databases/breast-cancer-wisconsin/>

[27] P.Wolfe, Finding the nearest point in a polytope, *Mathematical Programming*, Vol.11, 1976, pp.128-149.

### Appendix. Line Segment and Triangle Algorithms.

In this appendix we give two useful algorithms: one for computing the nearest point from the origin to a line segment; and one for computing the nearest point from the origin to a triangle.

Consider the problem of computing  $\tilde{z}$ , the nearest point of the line segment joining two points,  $z$  and  $\bar{z}$  from the origin. Expressing  $\tilde{z}$  in terms of a single real variable,  $\lambda$  as

$$\tilde{z} = \lambda\bar{z} + (1 - \lambda)z \quad (\text{A.1})$$

and minimizing  $\|\tilde{z}\|^2$  with respect to  $\lambda$ , it is easy to obtain the following expression for the optimal  $\lambda$ :

$$\lambda = \begin{cases} 1 & \text{if } \|\bar{z}\|^2 \leq \bar{z} \cdot z \\ \frac{\|z\|^2 - \bar{z} \cdot z}{\|\bar{z} - z\|^2} & \text{otherwise} \end{cases} \quad (\text{A.2})$$

It is also easy to derive an expression for the optimal value of  $\|\tilde{z}\|^2$ :

$$\|\tilde{z}\|^2 = \begin{cases} \|\bar{z}\|^2 & \text{if } \|\bar{z}\|^2 \leq \bar{z} \cdot z \\ \frac{\|z\|^2\|\bar{z}\|^2 - (\bar{z} \cdot z)^2}{\|\bar{z} - z\|^2} & \text{otherwise} \end{cases} \quad (\text{A.3})$$

Furthermore, suppose the following condition holds:

$$\bar{z} \cdot z < \|z\|^2 - \tilde{\epsilon} \quad (\text{A.4})$$

Let us consider two cases.

Case 1.  $\|\bar{z}\|^2 \leq \bar{z} \cdot z$ . In this case, we have, by (A.3) and (A.4),

$$\|z\|^2 - \|\tilde{z}\|^2 = \|z\|^2 - \|\bar{z}\|^2 \geq \|z\|^2 - \bar{z} \cdot z > \tilde{\epsilon} \quad (\text{A.5})$$

Case 2.  $\|\bar{z}\|^2 > \bar{z} \cdot z$ . Using (A.3) we get, after some algebra,

$$\|z\|^2 - \|\tilde{z}\|^2 = \frac{(\|z\|^2 - \bar{z} \cdot z)^2}{\|z - \bar{z}\|^2} \quad (\text{A.6})$$

Let  $\|z - \bar{z}\| \leq L$  for some finite  $L$ . (If  $z$  and  $\bar{z}$  are points of a polytope then such a bound exists. By (A.4) and (A.6) we get

$$\|z\|^2 - \|\tilde{z}\|^2 > \frac{\tilde{\epsilon}^2}{L^2} \quad (\text{A.7})$$

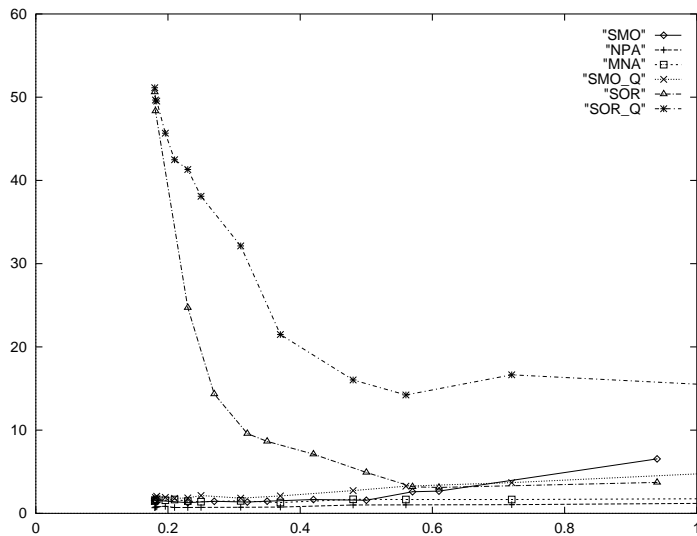


Figure 11: Wisconsin Breast Cancer data: Number of Kernel Evaluations  $\times 10^{-6}$  (vertical axis) shown as a function of  $M$  (horizontal axis)

Combining the two cases we get

$$\|z\|^2 - \|\tilde{z}\|^2 > \min\{\tilde{\epsilon}, \frac{\tilde{\epsilon}^2}{L^2}\} \quad (A.8)$$

Consider next, the problem of computing the nearest point of a triangle joining three points,  $P$ ,  $Q$  and  $R$  from the origin. We have found it to be numerically robust to compute the minimum distance from the origin to each of the edges of the triangle and the minimum distance from the origin to the interior of the triangle (if such a minimum exists) and then take the best of these four values. The first three distances can be computed using the line segment algorithm we described above. Let us now consider the interior. This is done by computing the minimum distance from the origin to the two dimensional affine space formed by  $P$ ,  $Q$  and  $R$ , and then testing whether the nearest point lies inside the triangle. Setting the gradient of  $\|P + (Q - P)\lambda_2 + (R - P)\lambda_3\|^2$  to 0 and solving for  $\lambda_1$  and  $\lambda_2$  yields

$$\lambda_2 = (e_{22}f_1 - e_{12}f_2)/den, \quad \lambda_3 = (-e_{12}f_1 + e_{11}f_2)/den \quad (A.9)$$

where:  $e_{11} = \|Q - P\|^2$ ,  $e_{22} = \|R - P\|^2$ ,  $e_{12} = (Q - P) \cdot (R - P)$ ,  $den = e_{11}e_{22} - e_{12}^2$ ,  $f_1 = (P - Q) \cdot P$ , and  $f_2 = (P - R) \cdot R$ . Let  $\lambda_1 = 1 - \lambda_2 - \lambda_3$ . Clearly, the minimum point to the affine space lies inside the triangle if and only if  $\lambda_1 > 0$ ,  $\lambda_2 > 0$  and  $\lambda_3 > 0$ ; in that case, the nearest point is given by  $\lambda_1 P + \lambda_2 Q + \lambda_3 R$  and the square of the minimum distance is  $P \cdot P - \lambda_2 f_1 - \lambda_3 f_2$ .

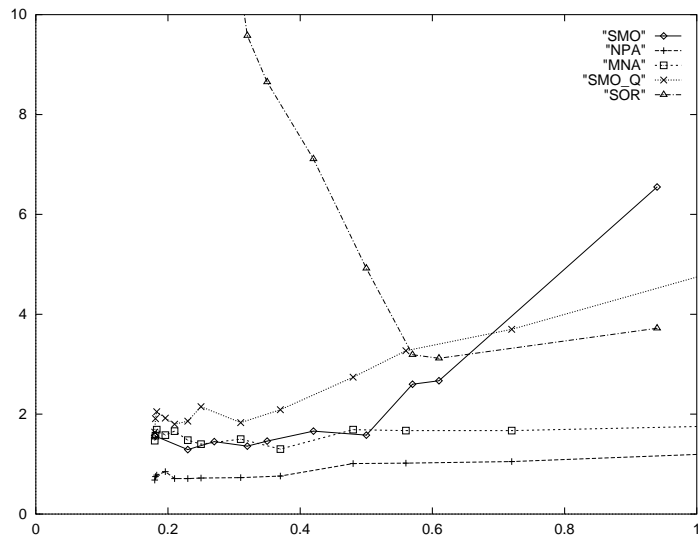


Figure 12: Figure 11 redrawn after leaving out SOR\_Q so as to compare other methods clearly.

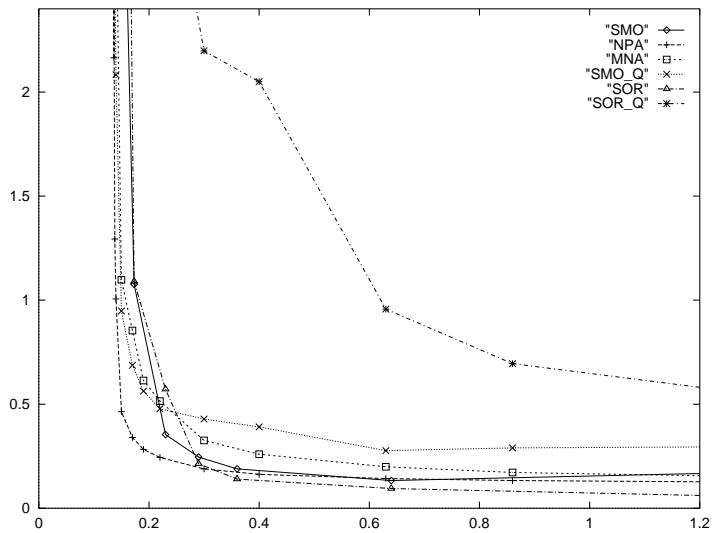


Figure 13: Two Spirals data: Number of Kernel Evaluations  $\times 10^{-6}$  (vertical axis) shown as a function of  $M$  (horizontal axis)

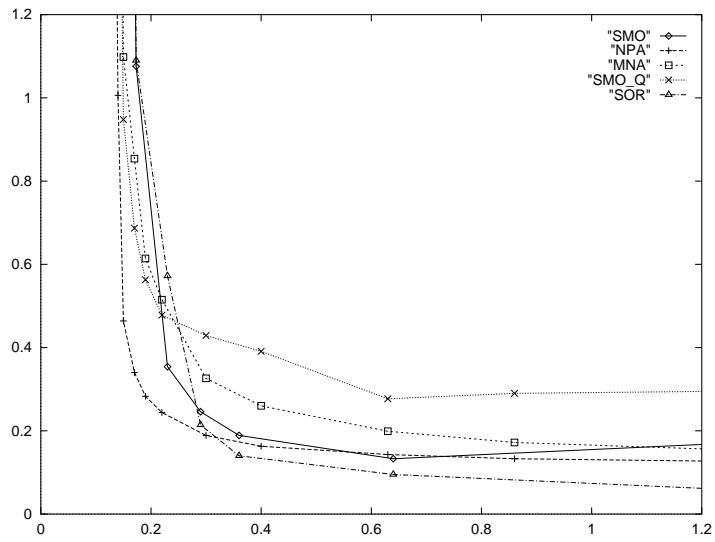


Figure 14: Figure 13 redrawn after leaving out SOR\_Q so as to compare other methods clearly.

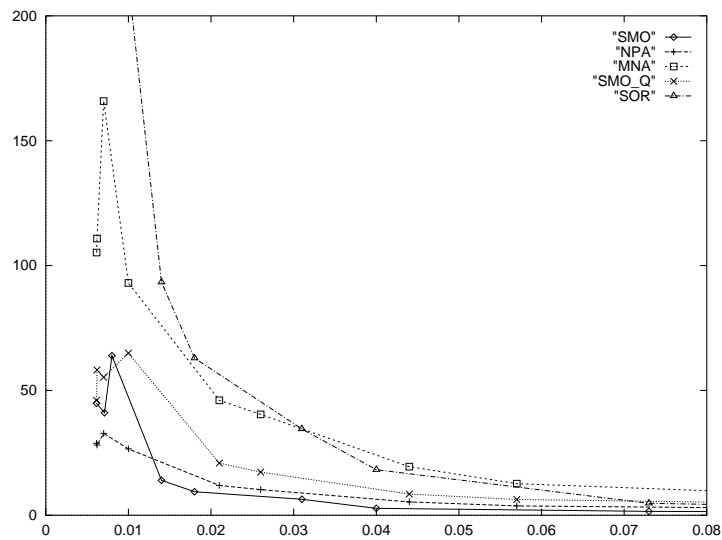


Figure 15: Checkers data: Number of Kernel Evaluations  $\times 10^{-6}$  (vertical axis) shown as a function of  $M$  (horizontal axis)

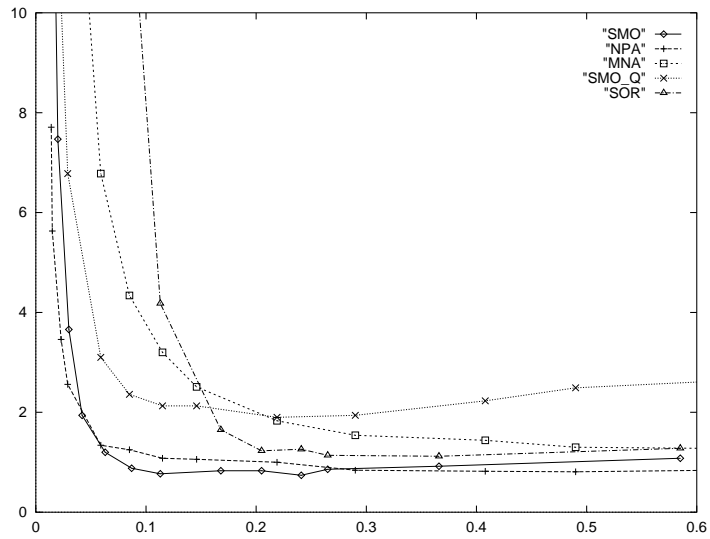


Figure 16: Adult-1 data: Number of Kernel Evaluations  $\times 10^{-7}$  (vertical axis) shown as a function of  $M$  (horizontal axis)

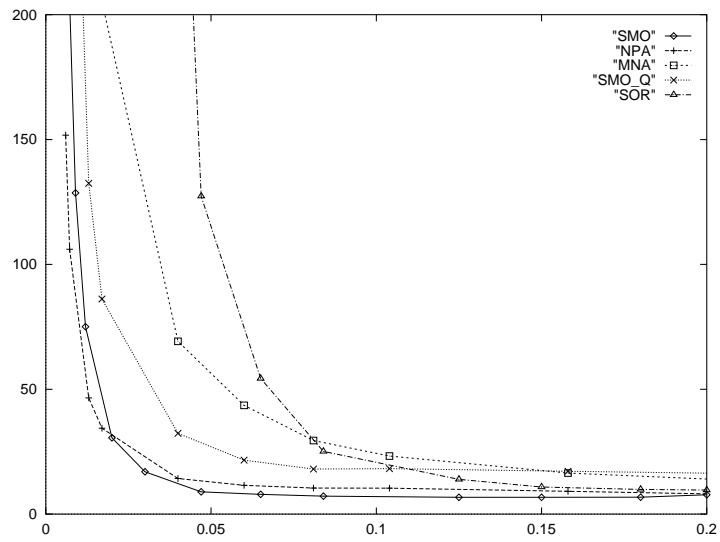


Figure 17: Adult-4 data: Number of Kernel Evaluations  $\times 10^{-7}$  (vertical axis) shown as a function of  $M$  (horizontal axis)

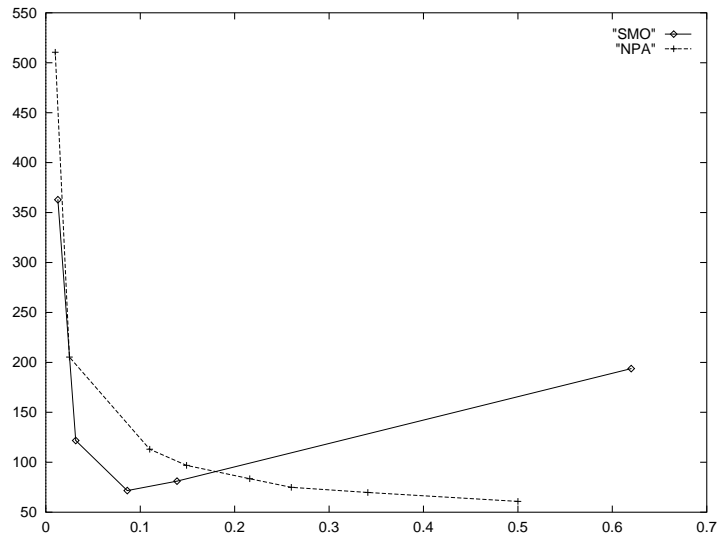


Figure 18: Adult-7 data: Number of Kernel Evaluations  $\times 10^{-7}$  (vertical axis) shown as a function of  $M$  (horizontal axis)

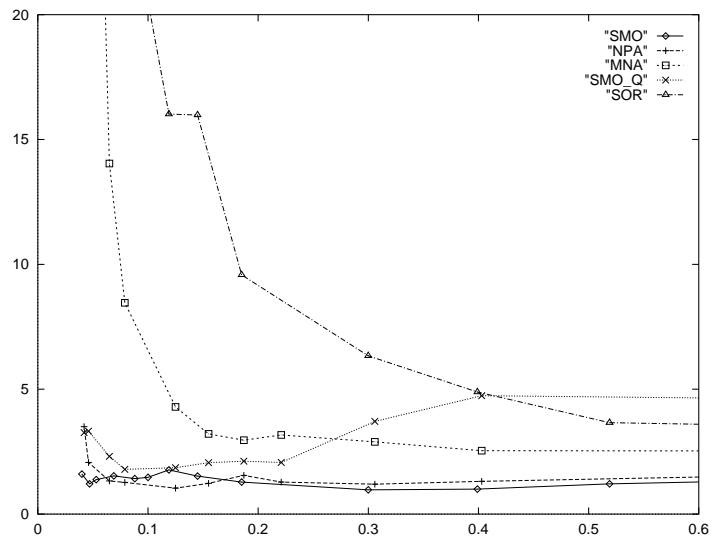


Figure 19: Web-1 data: Number of Kernel Evaluations  $\times 10^{-7}$  (vertical axis) shown as a function of  $M$  (horizontal axis)

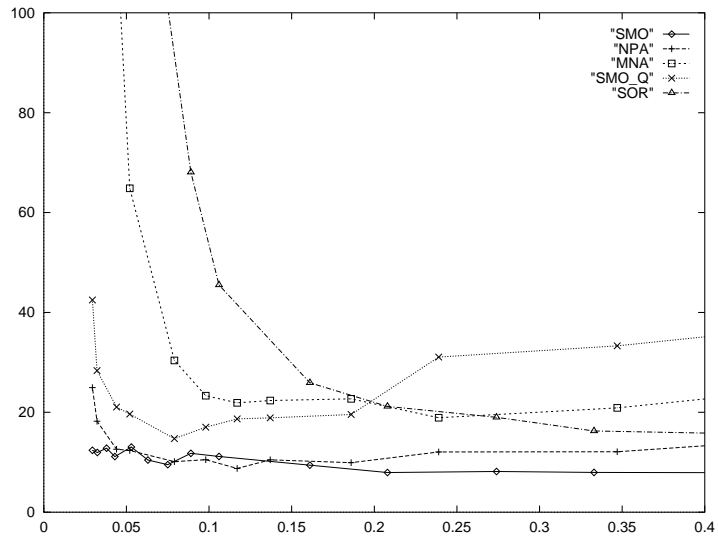


Figure 20: Web-4 data: Number of Kernel Evaluations  $\times 10^{-7}$  (vertical axis) shown as a function of  $M$  (horizontal axis)

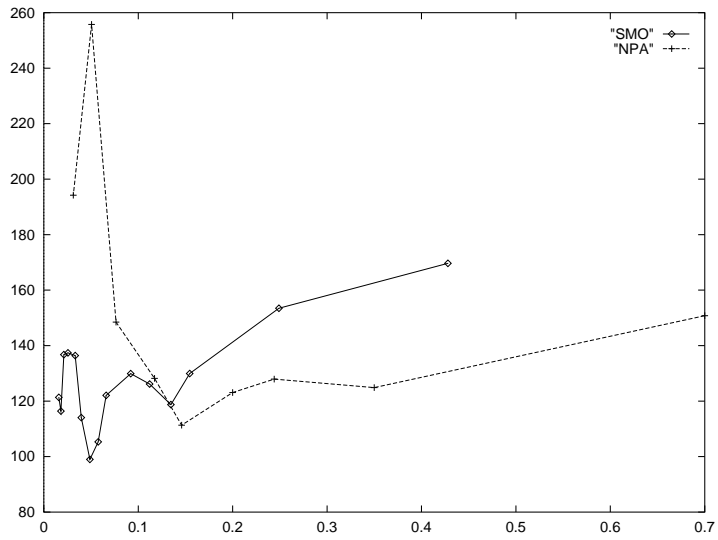


Figure 21: Web-7 data: Number of Kernel Evaluations  $\times 10^{-7}$  (vertical axis) shown as a function of  $M$  (horizontal axis)