

# Department of Electrical and Computer Systems Engineering

## Technical Report MECSE-29-2003

Increasing the step of the Newtonian Decomposition Method  
for Support Vector Machines

D.Lai, N.Mani and M.Palaniswami

**MONASH**  
UNIVERSITY

# Increasing the step of the Newtonian Decomposition Method for Support Vector Machines

\*D.Lai , \*N.Mani, +M.Palaniswami

\*Dept. of Electrical and Computer Systems Engineering  
Monash University, Clayton, Vic. 3168, Australia.

+Dept. of Electrical and Electronic Engineering,  
The University of Melbourne, Vic. 3010, Australia.  
{daniel.lai,n.mani@eng.monash.edu.au, swami@ee.mu.oz}

**Abstract:** The Newtonian method is a standard iterative method with quadratic convergence rates for solving large optimization problems. The Support Vector Machine problem formulation requires the solution of large datasets, which is ideal for the application of Newtonian methods. We point out that the Sequential Minimal Optimization (SMO) algorithm is a Newtonian Decomposition Method that is popular due to its efficiency but unfortunately does not scale too well with large data sets. The algorithm is an implementation of the decomposition method, which solves a sequence of sub problems instead of the entire problem at once. In this report, we introduce an extrapolation parameter to the SMO update method and investigate the effect on the rate of convergence of this algorithm. We first show that the SMO update method is Newtonian and that extrapolation ensures the update is norm reducing on the objective function. We also investigate the bounds of extrapolation and derive optimal estimates for this parameter. It was observed that choosing the working set pair according to some partial order does result in slightly faster speedups in algorithm performance.

## I. INTRODUCTION

The Support Vector Machines (SVM) developed by Vapnik[1] and co-workers has been shown to be a powerful supervised learning tool. The standard soft-margin Support Vector Machine is a binary classifier applied to classify a data set defined as,

$$\begin{aligned} \Theta &= \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2) \dots (\mathbf{x}_l, y_l)\} \\ \mathbf{x}_i &\in \mathbb{R}^n \\ y_i &= \{1, -1\} \end{aligned} \quad (1)$$

The SVM formulation is essentially a regularized minimization problem leading to the use of Lagrange Theory and quadratic programming techniques. The formulation defines a boundary separating two classes in the form of a linear hyperplane in data space where the distance between the boundaries of the two classes and the hyperplane is known as the margin of the hyperplane. This idea is further extended for data that is not linearly separable; where it is first mapped via a nonlinear function to a higher dimension feature space. Maximizing the margin of the hyperplane in either space is equivalent to maximizing the distance between the class boundaries. Vapnik[1] suggests that the form of the hyperplane,  $f(\mathbf{x}) \in F$  be chosen from family of functions with sufficient capacity. In particular,  $F$  contains functions for the linearly and non-linearly separable hyperplanes;

$$f(\mathbf{x}) = \sum_{i=1}^l w_i x_i + b \quad (2)$$

$$f(\mathbf{x}) = \sum_{i=1}^l w_i \phi(x_i) + b \quad (3)$$

The weight vector,  $\mathbf{w}$  in (3) is no longer the same expansion as in the linearly separable case (2). In fact, the non-linear mapping  $\phi: \mathbf{x} \subset \mathfrak{R}^n \rightarrow \mathfrak{R}^m$  and  $n, m \in [1, \infty)$  defines the mapping from data space to feature space. Hence the weights in feature space will have a one to one correspondence with the elements of  $\phi(\mathbf{x})$ . Now for separation in feature space, we would like to obtain the hyperplane with the following properties;

$$\begin{aligned}
 f(\mathbf{x}) &= \sum_{i=1}^l w_i \phi(\mathbf{x}_i) + b \\
 f(\mathbf{x}) &> 0 \quad \forall \quad i: y_i = +1 \\
 f(\mathbf{x}) &< 0 \quad \forall \quad i: y_i = -1
 \end{aligned} \tag{4}$$

The conditions in (4) can be described by a strict linear discriminant function, so that for each element pair in  $\Theta$ , we require that;

$$y_i \left( \sum_{i=1}^l w_i \phi(\mathbf{x}_i) + b \right) \geq 1 \tag{5}$$

The distance from the hyperplane to a support vector is  $\frac{1}{\|\mathbf{w}\|}$  and the distance between the support vectors of one class to the other class is simply  $\frac{2}{\|\mathbf{w}\|}$  by geometry. The soft-margin minimization problem relaxes the strict discriminant in (5) by introducing slack variables,  $\xi_i$  and is formulated as;

$$\text{P.1} \left\{ \begin{array}{l} \text{minimize } \mathfrak{J}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^l w_i^2 + C \sum_{i=1}^l \xi_i \\ \text{subject to } \left\{ \begin{array}{l} y_i \left( \sum_{i=1}^l w_i \phi(\mathbf{x}_i) + b \right) \geq 1 + \xi_i \\ \forall i = 1..l \end{array} \right. \end{array} \right. \tag{6}$$

We now apply Lagrange Theory to solve (6) giving us the Lagrange Primal problem of the following form;

$$\text{P.2} \left\{ \begin{array}{l} \text{minimize } \mathfrak{J}(\mathbf{w}, \alpha) = \frac{1}{2} \sum_{i=1}^l w_i^2 + \sum_{i=1}^l \alpha_i \left( y_i \left( \sum_{j=1}^l w_j \phi(x_j) + b \right) - 1 - \xi_i \right) + \sum_{i=1}^l \pi_i \xi_i \\ \text{subject to } \left\{ \begin{array}{l} \nabla \left( \frac{1}{2} \sum_{i=1}^l w_i^2 \right) + \nabla \left( \sum_{i=1}^l \alpha_i \left( y_i \left( \sum_{j=1}^l w_j \phi(x_j) + b \right) - 1 - \xi_i \right) \right) = 0 \\ \alpha_i, \pi_i \geq 0 \quad \forall i=1..l \end{array} \right. \end{array} \right. \tag{7}$$

The following dual optimization of a cost function written in terms of Lagrange Multipliers alone is usually implemented by incorporating the gradient condition into the cost function of P.2 and minimizing the following Lagrange dual in terms of Lagrange multipliers,  $\alpha$  alone;

$$\text{P.3} \left\{ \begin{array}{l} \min_{\alpha \in \mathbb{Q} \subset \mathbb{R}^l} \mathfrak{J}(\alpha) = \frac{1}{2} \sum_{i,j=1}^l \alpha_i \alpha_j y_i y_j K(x_i, x_j) - \sum_{i=1}^l \alpha_i \\ \text{where } \mathbb{Q} = \left\{ \alpha \mid 0 \leq \alpha_i \leq C, \sum_{i=1}^l \alpha_i y_i = 0, \forall i = 1..l \right\} \end{array} \right. \tag{8}$$

The separating hyperplane surface in (4) can now be written in terms of Lagrange Multipliers;

$$f(x) = \text{sign} \left( \sum_{i=1}^l \alpha_i y_i K(x, x_i) + b \right) \tag{9}$$

In this report, we point out that the SMO is in fact a Newton-method solving a different cost function than the problem originally stated in [2]. It was postulated by Chang [3] that update rules using the steepest gradient

direction might not give the optimal rate of convergence for a decomposition method. This leads us to investigate applying a larger Newtonian step controlled by an extrapolation parameter, which we will introduce in the following. We investigate applying our extrapolation parameter to the SMO update rule to increase the size of the step towards the minimum of the cost function. We also introduce the notion of *damping* and *extrapolation* regions to the SMO iteration by allowing our extrapolation parameter to vary and observe the resulting algorithm performance. We note here that several modern implementations of SMO and SVMlight such as [4] and LibSVM[5] are currently much faster than the original pseudocodes. However, the newer implementations run a basic combination of the original update rules using several advance heuristics and more efficient coding. Our work here investigates an improvement to the basic SMO update rule. We could apply this technique to modern implementations and investigate their convergences rates further.

The outline of this paper is as follows; in the next section, we review the Sequential Minimal Optimization algorithm and rewrite it in the familiar Newton update form. In Section 3, we try to justify the use of extrapolation to increase the speed of a general Newton update rule and further investigate its use with the bounded Newton form of SMO. The remaining sections will be left for our results and further discussions.

## II. SEQUENTIAL MINIMAL OPTIMIZATION (SMO): A NEWTONIAN DECOMPOSITION METHOD

We review the update rule used in SMO[2] which we show to be a hybrid element-wise Newton method in a decomposition algorithm setting. For the sake of simplicity, we will adopt the notations used by Platt in the original implementation for this summary. The working set size in SMO is fixed to two in order to enforce the equality constraint in (8) at every step of the iteration. We denote the working set as;

$$\alpha_w = \{\alpha_1, \alpha_2\} \quad \alpha_1, \alpha_2 \in \Re \quad (10)$$

Whenever we update the value of two multipliers, SMO maintains the equality constraint by ensuring;

$$\begin{aligned} &\alpha_1 y_1 + \alpha_2 y_2 = k \\ \text{where } &\begin{cases} \alpha_1 + \alpha_2 = k & \text{if } y_1 = y_2 \\ \alpha_1 - \alpha_2 = k & \text{if } y_1 \neq y_2 \end{cases} \end{aligned} \quad (11)$$

We make use of (11) to express the following for some iterative step,  $t$  where  $t > 0$ ;

$$\begin{aligned} &s = y_1 y_2 \\ &\alpha_1^{t+1} + s \alpha_2^{t+1} = \alpha_1^t + s \alpha_2^t = \gamma \end{aligned} \quad (12)$$

Platt defines the error of a training point as  $E_i$  where,

$$E_i = f(x_i) - y_i \quad \forall i = 1..l \quad (13)$$

The minimum of the cost function (8) subject to the equality constraint is first found in the direction of an arbitrary  $\alpha_2$  while the other multiplier  $\alpha_1$  is set to a value that maintains the equality constraint. We can derive the update rule for SMO by writing the cost function in terms of the working set members;

$$\mathfrak{J}(\alpha_1, \alpha_2) = -\alpha_1 - \alpha_2 + \frac{1}{2} K_{11} \alpha_1^2 + \frac{1}{2} K_{22} \alpha_2^2 + s K_{12} \alpha_1 \alpha_2 + y_1 \alpha_1 v_1 + y_2 \alpha_2 v_2 + \mathfrak{J}_{const} \quad (14)$$

$$\text{where } v_i = \sum_{j=3}^l y_j \alpha_j^t K_{ij} = f(x_i) - y_1 \alpha_1^t K_{1i} - y_2 \alpha_2^t K_{2i}$$

$$\mathfrak{J}_{const} = -\sum_{i=3}^l \alpha_i + \frac{1}{2} \sum_{i,j=3}^l \alpha_i \alpha_j y_i y_j K(x_i, x_j)$$

Then, the objective function expressed in  $\alpha_2$  entirely is,

$$\begin{aligned} \mathfrak{I}(\alpha_2) = & -\gamma + s\alpha_2 - \alpha_2 + \frac{1}{2}K_{11}(\gamma - s\alpha_2)^2 + \frac{1}{2}K_{22}\alpha_2^2 + sK_{12}(\gamma - s\alpha_2)\alpha_2 \\ & + y_1(\gamma - s\alpha_2)v_1 + y_2\alpha_2v_2 + \mathfrak{I}_{const} \end{aligned} \quad (15)$$

The minimum of (15) can be found with respect to  $\alpha_2$ , by taking the derivatives;

$$\nabla \mathfrak{I}(\alpha_2) = sK_{11}(\gamma - s\alpha_2) - K_{22}\alpha_2 + K_{12}\alpha_2 - sK_{12}(\gamma - s\alpha_2) - y_2v_2 + y_2v_1 - s + 1 = 0 \quad (16)$$

$$\nabla(\nabla \mathfrak{I}(\alpha_2)) = 2K_{12} - K_{22} - K_{11} \quad (17)$$

Expanding (16) and substituting for  $s$  and  $\gamma$ , we get

$$\alpha_2^{t+1}(K_{11} + K_{22} - 2K_{12}) = \alpha_2^t(K_{11} + K_{22} - 2K_{12}) + y_2(f(x_1) - f(x_2) + y_2 - y_1) \quad (18)$$

Finally, after some rearranging we obtain the update rule,

$$\alpha_2^{t+1} = \alpha_2^t + y_2 \frac{(E_1 - E_2)}{(K_{11} + K_{22} - 2K_{12})} \quad (19)$$

In order to ensure that  $\alpha_2 \in \mathbb{Q}$ , we enforce the following;

$$\alpha_2^{t+1} = \begin{cases} UB & \text{if } \alpha_2^{t+1} \geq UB \\ \alpha_2^{t+1} & \text{if } LB < \alpha_2^{t+1} < UB \\ LB & \alpha_2^{t+1} \leq LB \end{cases} \quad (20)$$

where;

$$\begin{aligned} UB &= \begin{cases} \min(C, \alpha_1 + \alpha_2) & \text{if } y_1 = y_2 \\ \min(C, C + \alpha_2 - \alpha_1) & \text{if } y_1 \neq y_2 \end{cases} \\ LB &= \begin{cases} \max(0, \alpha_1 + \alpha_2 - C) & \text{if } y_1 = y_2 \\ \max(0, \alpha_2 - \alpha_1) & \text{if } y_1 \neq y_2 \end{cases} \end{aligned} \quad (21)$$

The update for  $\alpha_1$  is then found by using (12);

$$\alpha_1^{t+1} = \alpha_1^t + s(\alpha_2^t - \alpha_2^{t+1, bounded}) \quad (22)$$

It can be seen here that the equality constraint is enforced at each step of the iteration through the use of (22). We now show that the update rule is Newtonian by first expanding (19);

$$\alpha_2^{t+1} = \alpha_2^t + y_2 \frac{(E_1 - E_2)}{(K_{11} + K_{22} - 2K_{12})} = \alpha_2^t + \frac{(y_2 f(x_1) - y_2 y_1 - y_2 f(x_2) + 1)}{(K_{11} + K_{22} - 2K_{12})} \quad (23)$$

Now, we simplify (16),

$$\begin{aligned} \nabla \mathfrak{I}(\alpha_2) &= sK_{11}(\gamma - s\alpha_2) - K_{22}\alpha_2 + K_{12}\alpha_2 - sK_{12}(\gamma - s\alpha_2) - y_2v_2 + y_2v_1 - s + 1 \\ &= sK_{11}(\alpha_1 + s\alpha_2 - s\alpha_2) - K_{22}\alpha_2 + K_{12}\alpha_2 - sK_{12}(\alpha_1 + s\alpha_2 - s\alpha_2) - y_2(f(x_2) + b - y_1\alpha_1K_{12} - y_2\alpha_2K_{22}) \\ &\quad + y_2(f(x_1) + b - y_1\alpha_1K_{11} - y_1\alpha_2K_{21}) - s + 1 \\ &= s\alpha_1K_{11} - K_{22}\alpha_2 + K_{12}\alpha_2 - s\alpha_1K_{12} - y_2f(x_2) + s\alpha_1K_{12} + \alpha_2K_{22} + y_2f(x_1) - s\alpha_1K_{11} + \alpha_2K_{21} - s + 1 \\ &= y_2f(x_1) - y_2f(x_2) - s + 1 \\ &= y_2f(x_1) - y_1y_2 - y_2f(x_2) + 1 \end{aligned}$$

Then using (17) we get the required relationship;

$$\begin{aligned}\alpha_2^{t+1} &= \alpha_2^t + \frac{(y_2 f(x_1) - y_2 y_1 - y_2 f(x_2) + 1)}{(K_{11} + K_{22} - 2K_{12})} \\ &\rightarrow \alpha_2^t - \frac{\nabla \mathfrak{Z}(\alpha_2)}{\nabla(\nabla \mathfrak{Z}(\alpha_2))}\end{aligned}\quad (24)$$

Since  $\alpha_2$  is arbitrary, we get the following relationship for  $\alpha_1$  ;

$$\begin{aligned}\nabla \mathfrak{Z}(\alpha_2) &= y_2 (E_1 - E_2) \\ &= -y_1 y_1 y_2 (E_2 - E_1) \\ &= -s y_1 (E_2 - E_1) \\ &= -s \nabla \mathfrak{Z}(\alpha_1)\end{aligned}\quad (25)$$

We note that (17) implies  $\nabla(\nabla \mathfrak{Z}(\alpha_2)) = \nabla(\nabla \mathfrak{Z}(\alpha_1))$  since  $\alpha_2$  is arbitrary. Then, we have from (22) the following result;

$$\begin{aligned}\alpha_1^{t+1} &= \alpha_1^t + s(\alpha_2^t - \alpha_2^{t+1, \text{bounded}}) \\ &= \alpha_1^t + s(\alpha_2^t - \alpha_2^t + \frac{\nabla \mathfrak{Z}(\alpha_2)}{\nabla(\nabla \mathfrak{Z}(\alpha_2))}) \\ &= \alpha_1^t + s \frac{\nabla \mathfrak{Z}(\alpha_2)}{\nabla(\nabla \mathfrak{Z}(\alpha_2))} \\ &= \alpha_1^t - s^2 \frac{\nabla \mathfrak{Z}(\alpha_1)}{\nabla(\nabla \mathfrak{Z}(\alpha_1))} = \alpha_1^t - \frac{\nabla \mathfrak{Z}(\alpha_1)}{\nabla(\nabla \mathfrak{Z}(\alpha_1))} \rightarrow \alpha_1^{t+1} \in \mathbb{Q}\end{aligned}\quad (26)$$

Then the update rule on the working set members can be written as;

$$\alpha_2^{t+1} = \alpha_2^t - \frac{\nabla \mathfrak{Z}(\alpha_2)}{\nabla(\nabla \mathfrak{Z}(\alpha_2))}\quad (27)$$

$$\alpha_1^{t+1} = \alpha_1^t - \frac{\nabla \mathfrak{Z}(\alpha_1)}{\nabla(\nabla \mathfrak{Z}(\alpha_1))}\quad (28)$$

We note from (8) that  $\alpha_1, \alpha_2 \in \mathbb{Q}$  implies that the multipliers are bounded at each iterative stage to ensure that they remain in the feasible region. Therefore, depending on the choice of  $\mathbb{Q}$ , the updates are optimal in the direction of the two multipliers on the feasible region. The standard Newton form[6] with the objective function  $F(\alpha)$  is given by;

$$\alpha^{t+1} = \alpha^t - \frac{F(\alpha^t)}{\nabla F(\alpha^t)}\quad (29)$$

By comparison, the forms of (27) and (28) seem to imply that from a Newtonian perspective, the objective function is instead the derivative of (8) by setting  $F(\alpha) = \nabla \mathfrak{Z}(\alpha)$ . This is now similar to the problem solved by Joachim's and others (see [5, 7]) where the derivative of the Lagrange Dual is solved to make computations easier and to avoid computing the quadratic function.

## III. EXTRAPOLATION: INCREASING THE STEP TOWARDS THE MINIMUM

## A. Extrapolation on the general Newton Method

We now introduce an extrapolation parameter  $\zeta$  to the general Newton iteration form below;

$$\alpha_i^{t+1} = \alpha_i^t - \zeta \frac{F(\alpha_i)}{\nabla F(\alpha_i)} \quad \forall i > 0 \quad (30)$$

We note here that when  $\zeta = 1$ , we recover the normal Newton step as in (29). We will show that the introduction of the extrapolation parameter to the general Newton method still ensures the iteration rule is norm reducing on the cost function. We then show that this extrapolation parameter is bounded to an interval when applied to a general Newton method. We first make use of a well-known result that we adapt in the following technical lemma;

Lemma III.1 (Ortega[8])

Let the function  $F : \mathbb{Q} \subset \mathfrak{R}^N \rightarrow \mathfrak{R}^M$  be Gateux-differentiable (G-differentiable) on the convex set  $\mathbb{Q} \in \mathfrak{R}^N$ . Then for any vector  $x, y \in \mathbb{Q}$ , we have the following;

$$\|Fy - Fx\| \leq \sup_{0 \leq t \leq 1} \|\nabla F(x + t(y-x))\| \|y-x\| \quad (31)$$

Proof:

First, assume that the supremum is bounded from above, i.e.

$$S = \sup_{0 \leq t \leq 1} \|\nabla F(x + t(y-x))\| < \infty$$

Then for an arbitrary  $\varepsilon > 0$ , we can choose a  $t \in [0, 1]$  such that

$$\begin{aligned} \|F(x + t(y-x)) - Fx - t\nabla F(x + t(y-x))(y-x)\| &\leq \varepsilon t \|y-x\| \\ \|F(x + t(y-x)) - Fx\| &\leq \sup_{0 \leq t \leq 1} t \|\nabla F(x + t(y-x))(y-x)\| + \varepsilon t \|y-x\| \\ &\leq St \|y-x\| + \varepsilon t \|y-x\| \end{aligned}$$

Since F is G-differentiable on a convex set, it is also continuous in t so that for any  $\gamma$  the following holds,

$$\|F(x + \gamma(y-x)) - Fx\| \leq S\gamma \|y-x\| + \varepsilon\gamma \|y-x\| \quad (32)$$

When  $\gamma = 1$ , we recover (31)

$$\begin{aligned} \|Fy - Fx\| &\leq S\|y-x\| + \varepsilon\|y-x\| \\ &\leq \sup_{0 \leq t \leq 1} \|\nabla F(x + t(y-x))\| \|y-x\| + \varepsilon\|y-x\| \\ &\leq \sup_{0 \leq t \leq 1} \|\nabla F(x + t(y-x))\| \|y-x\| \end{aligned}$$

Assume  $\gamma < 1$  and since  $\nabla F$  exists everywhere on the convex set, there is a  $\lambda$  such that  $\lambda \in (\gamma, 1)$  and we have;

$$\begin{aligned} \|F(x + \lambda(y-x)) - F(x + \gamma(y-x)) - \nabla F(x + \lambda(y-x))(\lambda - \gamma)(y-x)\| &\leq \varepsilon(\lambda - \gamma)\|y-x\| \\ \rightarrow \|F(x + \lambda(y-x)) - F(x + \gamma(y-x))\| &\leq (\lambda - \gamma) \|\nabla F(x + \lambda(y-x))\| \|y-x\| + \varepsilon(\lambda - \gamma)\|y-x\| \\ &\leq (\|\nabla F(x + \lambda(y-x))\| + \varepsilon)(\lambda - \gamma)\|y-x\| \end{aligned}$$

But from (32), we have;

$$\|F(x + \lambda(y-x)) - Fx\| \leq (\|\nabla F(x + \lambda(y-x))\| + \varepsilon)\gamma\|y-x\|$$

This implies that  $\lambda > \gamma > 1$  which contradicts our previous assumption of  $\lambda, \gamma$  and so (31) must hold  $\square$

The proof of Lemma III.1 makes use of the definition of G-differentiability of functions on a convex set, which is defined in the Appendix. We now show that extrapolation ensures that the update rule is norm reducing in the following theorem.

Theorem III.I

Let the function  $F : \mathbb{Q} \subset \mathfrak{R}^N \rightarrow \mathfrak{R}^N$  be G-differentiable at a point  $x$  and  $\nabla F(x)$  be invertible. Then there exists a

scalar  $\zeta > 0$  such that for any vector  $x \in \mathbb{Q}$  and  $y = \frac{F(x)}{\nabla F(x)}$ ;

$$\|F(x - \zeta y)\| \leq \|Fx\| \quad (33)$$

Proof:

Since F is G-differentiable, we have the following relationship;

$$F(x - t\zeta y) - Fx + \psi t\zeta y \leq \nabla Fx \quad (34)$$

Here  $\psi \in L(\mathfrak{R}^n, \mathfrak{R}^m)$  is some linear operator and we get equality only at  $t=0$ . We now define an arbitrary norm on  $\mathfrak{R}^N$  and rearranging (34) obtain the relationship;

$$\|F(x - t\zeta y) - Fx - \nabla Fx\| \leq -t\zeta\|\psi y\|$$

Let  $t=1$  and we now have for some  $\varepsilon > 0$ ;

$$\begin{aligned} \|F(x - \zeta y) - Fx\| &\leq \|F(x - \zeta y)\| - \|Fx\| \leq -\varepsilon\zeta\|\psi y\| \\ \rightarrow \|F(x - \zeta y)\| &\leq \|Fx\| - \varepsilon\zeta\|\psi y\| \end{aligned}$$

Since  $\varepsilon$  is an arbitrary positive value, clearly (33) holds if and only if  $\zeta > 0$ .  $\square$

However, extrapolation is only norm-reducing provided that the parameter is positive which gives us a lower bound on the range of possible values. We now compute the upper bound for this parameter, by showing in the following theorem that the norm reduction property holds only if  $\zeta$  lies in a bounded interval.

Theorem III.II

Let the function  $F : \mathbb{Q} \subset \mathfrak{R}^N \rightarrow \mathfrak{R}^N$  be G-differentiable at the points  $x, y \in \mathbb{Q}$  and suppose the following holds,

$$\|\nabla F(x - ty) - \nabla F(x)\| \leq \psi t\|y\| \quad \forall t \in [0,1] \quad (35)$$

then (33) holds for all  $\zeta \in \left[0, \frac{\|Fx\|}{\psi\|y\|^2}\right]$ .

Proof:

Since F is G-differentiable, then (31) in Lemma III.1 gives the following for any  $t \in [0,1]$ ;



$$\|F(x - \zeta y) - Fx\| \leq \sup_{0 \leq t \leq 1} t \|\nabla F(x - t(x - \zeta y))\| \|x - \zeta y\| \quad (36)$$

We have;

$$\begin{aligned} \sup_{0 \leq t \leq 1} t \|\nabla F(x - t(x - \zeta y))\| \|x - \zeta y\| &= \sup_{0 \leq t \leq 1} t \|\nabla F(x - t(x - \zeta y))(x - \zeta y)\| \\ &= \sup_{0 \leq t \leq 1} t \|x \nabla F(x - t(x - \zeta y)) - \zeta y \nabla F(x - t(x - \zeta y))\| \\ &\leq \sup_{0 \leq t \leq 1} t \|\zeta y \nabla F(x - t(x - \zeta y))\| \\ &\leq \sup_{0 \leq t \leq 1} t \zeta \|y\| \|\nabla F(x - \zeta ty) - \nabla F(tx)\| \end{aligned}$$

Then from (36) and noting that  $\zeta > 0$ , we have

$$\begin{aligned} \|F(x - \zeta y) - Fx\| &\leq \sup_{0 \leq t \leq 1} t \zeta \|y\| \|\nabla F(x - \zeta ty) - \nabla F(tx)\| \\ &\leq \sup_{0 \leq t \leq 1} t \zeta \|y\| \|\nabla F(x - ty) - \nabla F(tx)\| \end{aligned} \quad (37)$$

Let  $t = 1$ , then if (33) holds, we have the following upper bound;

$$\begin{aligned} \|F(x - \zeta y) - Fx\| &\leq \zeta \|y\| \|\nabla F(x - ty) - \nabla F(x)\| \\ &\leq \zeta \psi \|y\|^2 \leq \|Fx\| \\ \rightarrow \zeta &\leq \frac{\|Fx\|}{\psi \|y\|^2} \end{aligned}$$

Thus for Theorem III.I to hold  $\zeta$  is constrained to the interval  $\left[0, \frac{\|Fx\|}{\psi \|y\|^2}\right]$ . We define the linear operator  $\psi$  to be identity matrix I when using the Euclidean norm, so that the element wise update rule is simplified. The bounds for the extrapolation parameter  $\zeta$  can now simplified to  $0 < \zeta < \frac{\nabla^2 Fx}{Fx}$   $\square$

The norm reducing theorems we derived above are in terms of minimization of the cost function and the domain space where the update rule is a mapping on the entire domain. It is trivial to show that they hold verbatim for maximizations of the cost function by simply replacing  $F(x)$  with its negative and rewriting them to show that extrapolation is also norm maximizing. The theorems above show that extrapolation can increase the step to the minimum of the cost function for a Newtonian method by ensuring the update is norm reducing. In the following, we will test our proposed extrapolation on the standard SMO algorithm and measure the improvement in performance over some benchmark datasets. The extrapolated SMO update rule for the updates of the working set is now just;

$$\alpha_w^{t+1} = \begin{pmatrix} \alpha_1^{t+1} \\ \alpha_2^{t+1} \end{pmatrix} = \begin{pmatrix} \alpha_1^t \\ \alpha_2^t \end{pmatrix} - \begin{pmatrix} \frac{\zeta y_1 (E_2 - E_1)}{\eta} \\ \frac{\zeta y_2 (E_1 - E_2)}{\eta} \end{pmatrix} \quad (38)$$

#### IV. EXPERIMENTAL METHOD

In the series of experiments, a Pentium IV, 1.5 GHz computer with 256MB RAM was used. Our proposed idea was implemented in the form of an extrapolated SMO on Visual C++ 6.0 with kernel caching. We recorded the number of iterations and epochs required till convergence was established. Here iterations refer to the number of successful updates on the Lagrange multipliers and epochs refer to the number of loops which the algorithm goes

through in the main program. We have kept our algorithm consistent with the original SMO in order to test the effect of our proposed extrapolation on the update rule. Thus, we need only to compare the number of iterations and epochs to determine the effectiveness of our extrapolated SMO. The CPU times are different from Platt's report due to the different architecture, which this experiment is executed on. We ran our proposed SMO on the UCI adult data set [9] and fish data set with a tolerance of 0.001. All results have been averaged out over 10 runs in order to minimize the effects of random ordering imposed by SMO heuristics on the speed of the algorithm.

For the UCI datasets we set  $C=1$  and applied a Gaussian kernel for separation. We investigated the algorithm performance over the range  $0.8 < \zeta < 1.8$  and report the results for the first 3 datasets in Fig 1. We also observed the behavior of the objective function (8) comparing the effect of  $\zeta$  on the number of epochs required till algorithm termination. We then ran a series of experiments of UCI adult set 1-5 at a higher tolerance to investigate the effect of extrapolation on the sensitivity of the algorithm to violating points. Next, we report the results for training a SVM to classify the UCI web 1 and web 2 ( Fig 5) data using a Gaussian kernel with  $C=5$ .

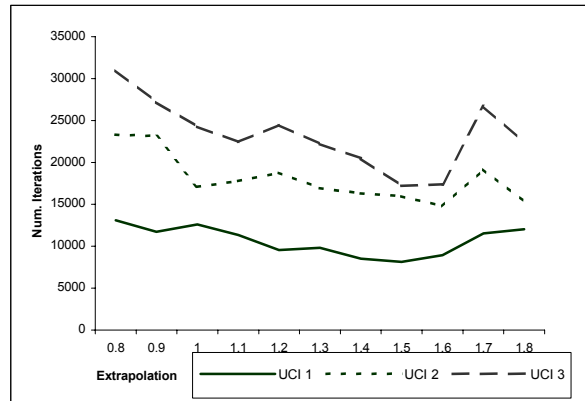


Fig 1: Number of iterations required for SMO across 3 datasets. Damped( $\zeta < 1$ ) and extrapolation( $\zeta > 1$ ) regions shown here.

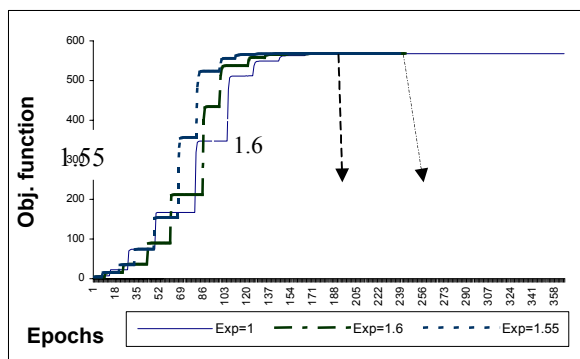


Fig 2: Behaviour of the Wolfe objective function for classical SMO ( $\zeta = 1$ ) and the extrapolated SMO. Arrows indicate the points of termination for the algorithm and the respective extrapolation parameters.

Extrapolation	UCI 1	UCI 2	UCI 3	UCI 4	UCI 5
0.8	13112	22453	30980	52851	76868
0.9	11726	22557	27208	46656	62693
1	12621	21455	24258	40409	72384
1.1	11340	17657	22439	37320	60704
1.2	9567	19372	24455	38859	57713
1.3	9818	15100	22183	35841	47314
1.4	8514	17973	20462	35256	55683
1.5	8151	15700	17209	32405	54063
1.6	8944	17070	17400	33391	53454
1.7	11537	18073	26704	36302	47912
1.8	12037	15884	22370	38567	53838

Table 1: Number of iterations till convergence using extrapolation for UCI adult datasets 1-5 at a tolerance of  $10^{-4}$

We then examined the effect of arranging or sorting the order of support vectors in the cache to be chosen for update. This implicitly imposes some form of partial ordering on the working set. We sorted the Support Vectors in order of increasing and decreasing magnitude and compare the resulting algorithm speed against the random choice of maximal violators. The results are shown for number of iterations in Fig 3 while Fig 4 shows the best improvement in speed when sorting the Support Vectors in decreasing order. We then tested the speed of the algorithm across a range of C for an arbitrary UCI dataset; in this case we show the results for UCI 3 in Fig 6.

We also ran our fish dataset across a range of kernels since it was highly inseparable. We have found that  $C=1$  produced an acceptable accuracy on the test set [10] when we tested this proposed algorithm using the Gaussian and Linear Kernels. For the Gaussian kernel, when  $\sigma=10$ , we discovered that  $\zeta$  lay in the range of  $0.5 < \zeta < 2.0$  and  $\zeta > 2.0$  caused the iteration to oscillate and not converge. Furthermore, we obtained a large improvement of about 50% less iterations required using extrapolated SMO against the normal SMO for this kernel Fig 7. It was noted by Lin[11] that SMO may be slow for linear kernels from the results of his co-workers. We thus applied a linear kernel to the data and observed in Fig 8 an improvement of about 25% in the speed of the algorithm.

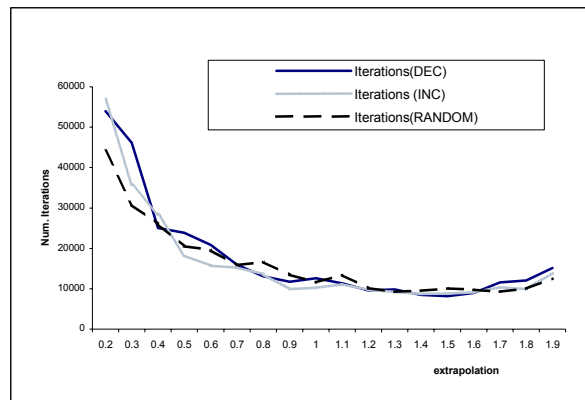


Fig 3: Number of successful updates for extrapolated and damped SMO using random, increasing, decreasing cache sorting on UCI 1.

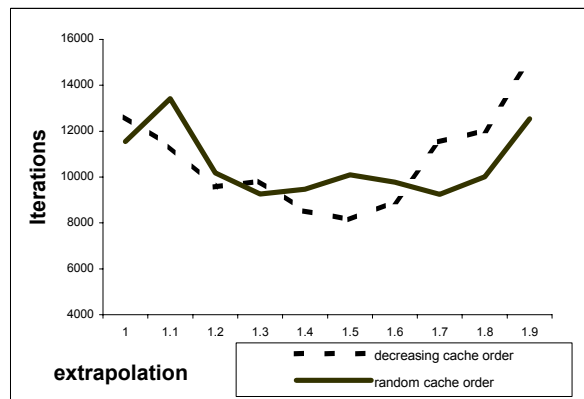


Fig 4: Performance of SMO in extrapolation on UCI 1 for two cache modes.

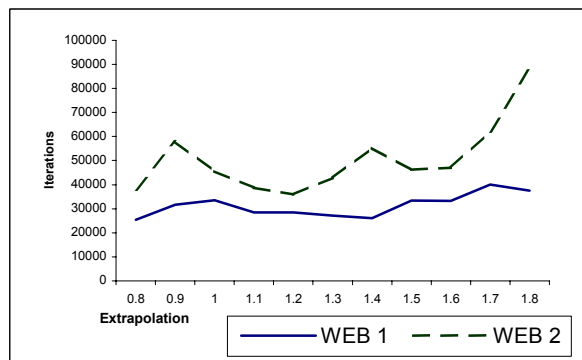


Fig 5: Extrapolated SMO on UCI web1 ( 2477 examples ) and web2(3470 examples)

## V. DISCUSSION

Platt inadvertently proposed a decomposition algorithm with a Newton iterate rule and solved the problem of applying a Newton rule to large datasets without the memory constraint problem. However, we believe that the rates of convergence of SMO are not quadratic like Newton methods due to the decomposition aspect of the problem. In fact, a worse case convergence rate has been shown to be linear[12]. We have shown that the update rule of SMO is simply just a bounded form of Newtonian iteration with clever heuristics besides being related to Bergman methods as claimed by [2]. From Fig 3, we can observe that choosing values of  $\zeta < 1$  results in a *damped* SMO where the update rule is not optimal in a Newtonian sense. Empirical results show that values of  $\zeta < 1$  (larger damping constant) not too far away from unity do give faster rates of convergence. However, if  $\zeta$  is too small, the rate of convergence becomes very much slower and is no longer practical. We say the algorithm is over-damped and convergence is very slow. We found that when testing on UCI datasets, at  $\zeta = 0.1$  the algorithm took excessively long to converge and for smaller values it had not converged after 5 hours.

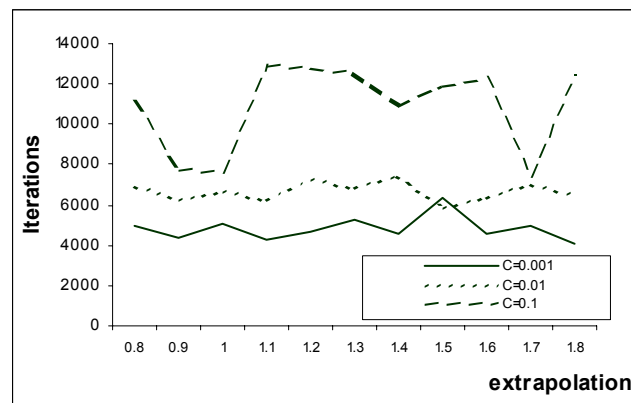


Fig 6: Extrapolated-SMO against normal SMO( $\zeta=1$ ) for UCI 3 data set over a range of C parameters.

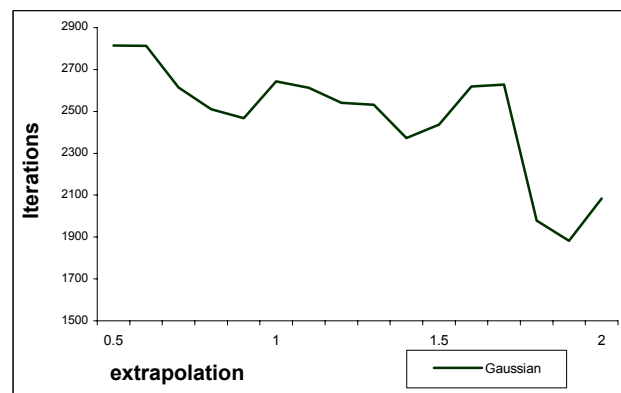


Fig 7: Effect of extrapolation on fish data set for Gaussian kernel with  $\sigma=10$

Overall, values of  $\zeta > 1$  give better improvement to the speed of the algorithm. In this region, the algorithm is in *extrapolation* and as long as  $\zeta$  satisfies the upper bound the algorithm converges. We found that for most cases  $\zeta \geq 2$  resulted in oscillation and non-convergence. In the case of the fish dataset, the upper bound for  $\zeta$  was around 1.95 for a Gaussian kernel with a kernel width of 10, before the algorithm stopped converging. On closer inspection, we found that a kernel width of this size still resulted in a highly inseparable feature space with many Lagrange multipliers at bound. This leads us to believe that extrapolation works better when the optimal solution has more variables at the boundary. However, we have a problem of quantifying this due to the random nature of the working set choice by Platt.

The introduction of the cache to store Support Vectors presents this randomness, which affects the working set choice, and hence also our empirical results. In the original SMO algorithm, Platt stored the errors of the training examples that had Lagrange Multipliers, which were Support Vectors. In our algorithm, our cache imposes an ordering on the Support Vectors obtained at each epoch of the iteration. We observed from Fig 4 empirically that sorting the Support Vectors to be updated according to some criteria in SMO does vary the

speed as compared to a random selection of Support Vectors. This result requires us to do several runs to average out the effect imposed by sorting in the cache. We then ran a small experiment to remove the effect of this random choice of working sets. This was done on the UCI adult 1 dataset over a range of  $C$  and the results are reported in Fig 9-Fig 11. It can be seen that extrapolation works better for larger values of  $C$ . Geometrically, as  $C$  increases the size of the feasible region increases while the Newtonian step remains the same. Hence during the initial stages of the iteration, we would like to determine the variables at bound quickly and the addition of the extrapolation parameter seems to boost the variables to the bound thus increasing the rate of convergence. Further empirical evidence can also be seen by reexamining Fig 6 where extrapolation does not work too well for small values of  $C$ . In fact, our view of extrapolation seems to be as a boost to the basic Newtonian step similar to our previous work on linear stationary iterative methods [10].

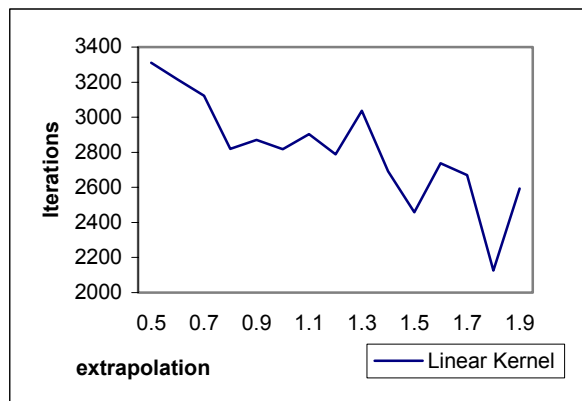


Fig 8: Effect of extrapolation on fish data set for linear kernel

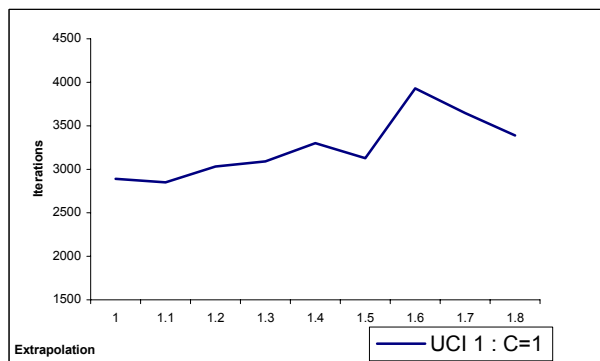


Fig 9: Extrapolation on UCI 1 without random heuristics for Gaussian kernel with C=1

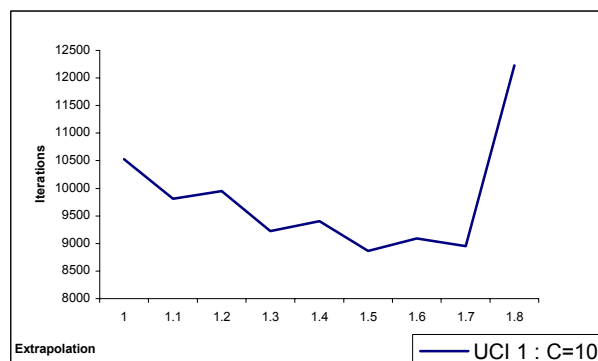


Fig 10: Extrapolation on UCI 1 without random heuristics for Gaussian kernel with C=10

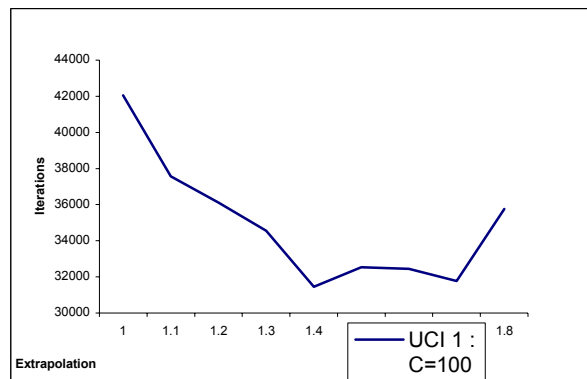


Fig 11: Extrapolation on UCI 1 without random heuristics for Gaussian kernel with C=100

The results from the experiments on the sorting indicate that the sequence in which the variables are updated do affect the rates of convergence of the decomposition algorithm. Certain orders give faster convergence while others slow the algorithm down tremendously. However, unfortunately we have yet to understand this observation clearly from a mathematical viewpoint or provide concrete justification for the optimal choice of sorting order. We do know however that, if Platt used a fixed size cache for errors, not all the Support Vectors were present in the cache at any one time, so indirectly the choice of points available for update was random and restricted to the size of the cache. In our case, we have the full array of Support Vectors to choose from and this presents some form of partial ordering on the domain of Lagrange Multipliers. We believe the analysis of the dynamics of the SMO algorithm under the choice of updating pairs would be worthwhile and more involving in order to increase its performance.

The analysis of the optimal value of extrapolation parameter  $\zeta$ , is an ongoing research topic for us. In all our benchmark tests, the use of  $\zeta$  results in an average improvement of about 35% compared to the original SMO without extrapolation. Intuitively, the largest value of  $\zeta$  is 2, which geometrically means double the step size and anything larger would result in us not taking a step to the minimum of the objective function. Finer values of  $\zeta$  were observed to give even better algorithm performance. This observation underscores the need for a theory to define the optimal value of  $\zeta$  and relate it in some way to the kernel matrix. We have given the necessary condition for  $\zeta$  to ensure that the iteration strictly reduces the objective function, however we do not know for sure yet whether this condition is sufficient to ensure the convergence of the algorithm. We believe however that the decomposition method is convergent as long as the basic update rule used is norm reducing on the objective function. Lin [11] has given a convergence proof for [7] in terms of a counting method, but we intend to relate the norm-reduction properties of the iteration to the convergence of the decomposition method.

It should be noted here that a consequence of the above observation is that we could use a different update rule instead of the Newton rule and apply the decomposition method to it. We can then investigate the optimal update rule in the decomposition setting with the dynamic determination of the threshold  $b$ . i.e. always maintaining the validity of the equality constraint or apply it to formulations that do not require explicit preservation of the equality constraint e.g.[10]. We believe that different working sets require different degrees or values of extrapolation hence the possibility of an adaptive extrapolation algorithm. The problem with this is that computing the optimal extrapolation parameter at each iterative step may add unnecessary overhead to the iteration and degrade the algorithm speed instead. We plan to look into this possibility in the near future.

## VI. CONCLUSION

We have proposed applying extrapolation and observed some improvement to the rates of convergence with better improvements for difficult problems i.e. larger values of  $C$ . We believe that the speed of this algorithm can be further improved if we apply newer methods that improve the computation of the threshold in SMO. Our future work will concentrate on investigating the proper choice of working set pairs and the theoretical rates of convergence using an extrapolated Newtonian update rule. In particular, we will investigate the relationship between optimal values of the extrapolation parameter and the structure of the feasible region formed by the constraint set.

## VII. APPENDIX

## A. Gateux-differentiable functions

Definition VII.1 (Ortega[8])

A mapping  $F : D \subset R^n \rightarrow R^m$  is G-differentiable at an interior point of D if there is some linear operator,  $A \in L(R^n, R^m)$  so that for any  $h \in R^n$ ,

$$\lim_{t \rightarrow 0} \frac{1}{t} (F(x + th) - Fx - tAh) = 0 \quad \bullet$$

## REFERENCES

- [1] V. N. Vapnik, *The nature of statistical learning theory*, 2nd ed. New York: Springer, 2000.
- [2] J. Platt, "Fast Training of Support Vector Machines Using Sequential Minimal Optimization," in *Advances in Kernel Methods-Support Vector Learning*, B. Schölkopf, C. J. C. Burges, and A. J. Smola, Eds.: Cambridge MIT Press, 1998, pp. 185-208.
- [3] C.-C. Chang, C.-W. Hsu, and C.-J. Lin, "The analysis of decomposition methods for support vector machines," *Neural Networks, IEEE Transactions on*, vol. 11, pp. 1003-1008, 2000.
- [4] S. K. S. S. Keerthi, C. Bhattacharyya and K.R.K. Murthy, "Improvements to Platt's SMO algorithm for SVM classifier design,," Control Division, Dept. of Mechanical Engineering, National University of Singapore CD-99-14, 1999.
- [5] C.-J. Lin, "LIBSVM," <http://www.csie.ntu.edu.tw/~cjlin/>.
- [6] L. A. Hageman and D. M. Young, *Applied iterative methods*. New York: Academic Press, 1981.
- [7] T. Joachims, "Making Large Scale Support Vector Machine Learning Practical," in *Advances in Kernel Methods - Support Vector Learning*, B. Schölkopf, C. J. C. Burges, and A. J. Smola, Eds.: Cambridge , MIT Press, 1998, pp. 169-184.
- [8] J. M. Ortega and W. C. Rheinboldt, *Iterative solution of nonlinear equations in several variables*. New York,: Academic Press, 1970.
- [9] C. L. M. Blake, C.J., "UCI Repository of machine learning databases Irvine, CA," University of California, Department of Information and Computer Science, 1998.
- [10] D. Lai, M. Palaniswami, and N. Mani, "Fast linear stationary methods for automatically biased support vector machines," presented at Neural Networks, 2003. Proceedings of the International Joint Conference on, 2003.
- [11] C.-J. Lin, "On the convergence of the decomposition method for support vector machines," *Neural Networks, IEEE Transactions on*, vol. 12, pp. 1288-1298, 2001.
- [12] C.-J. Lin, "Linear convergence for a decomposition method for Support Vector Machines," November 2001.