

# Department of Electrical and Computer Systems Engineering

## Technical Report MECSE-7-2003

Matrix Formulation for the Support Vector Machine  
Classifier

D. Lai and N. Mani

**MONASH**  
UNIVERSITY

# Matrix Formulation for the Support Vector Machine Classifier

D. LAI, N.MANI

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

23/7/2003

**Abstract:** In this paper, we investigate solving the constrained quadratic program for Support Vector Machine classification formulation by solving a constrained set of linear equations written in matrix form. We have applied this form in our previous work and have observed its validity through empirical results. Several other researchers have proposed optimization algorithms, which solved a linear minimization form similar in nature. We first derive this linear form through the use of convex mathematical programming and rewrite it in matrix form. We attempt to reconcile our form with that which has been popularly used.

## I. INTRODUCTION

The Support Vector Machine(SVM) formulation is a supervised learning formulation introduced by Vapnik[1] for pattern recognition. This formulation embodies the Structural Risk Minimization principle proposed by Vapnik in his pioneering of statistical learning theory where the performance of the SVM binary classifier is determined by controlling its capacity and minimizing the training error on the data set. The SVM classifier problem is generally treated using mathematical programming methods which results in a constrained quadratic program obtained from Lagrange Theory. The constrained Lagrange for the SVM classifier is a quadratic program that generally requires the use of interior point methods, activesets or some form of chunking[2] which are somewhat complex to implement for first time users of Support Vector Machines. Furthermore, the standard Lagrange Dual results in an objective function subject to an equality constraint and bounded Lagrange Multipliers which require enforcement at each stage of the iteration.

Several researchers have proposed optimization algorithms[3-5] that solve a linear program instead. They have partitioned their data into subsets and iterated on individual subsets. This amounts to solving a decomposed problem, hence the apt name of decomposition methods. The convergence of the decomposition methods have been examined by[6]. For now, this method seems to be popular and easy to implement with comparable optimization times.

In this paper, we proceed further by showing that this linear program can be written as a set of linear equations which surprisingly results in minimizing the training error on a data set. Furthermore we show this is equivalent to obtaining the solution to SVM classification which also enforces Structural Risk Minimization. We first attempt to explore the duality properties of the SVM quadratic program. We then show that solving a constrained minimization of training errors directly satisfies the optimal solution requirements for the dual Lagrange program. The resulting minimization program is a linear program with implicit maximization of the margin of the hyperplane. Our linear program constitutes solving a set of linear equations which can be written in matrix form and allows simple iterative algorithms to obtain the solution. We give a geometrical interpretation of the feasible region of solutions resulting from the mathematical programs in the hope of providing a better understanding of the underlying mechanics of optimization algorithms. The geometrical view is specific to the set of solutions and is different from previous work [7] which examines the construction of the classifier based on the geometric distribution of data.

This paper is divided in to sections as follows; section 2 will review the Support Vector Machine classifier from a mathematical programming perspective and draw attention to the geometric properties resulting from the constraints on the resulting objective functions. Section 3 is devoted to the derivation of the linear program and the matrix form while Section 4 will investigate the solution set further. The

remainder of the paper will be devoted to a discussion on the possibilities of the matrix form and future directions for research.

## II. SUPPORT VECTOR MACHINE CLASSIFICATION

### A. Overview of the Support Vector Classification Formulation

In binary classification, we are given a set of training data which have been labelled according to two classes. For historical reasons and the ease of modelling, the labels are chosen to be +1 and -1. The data to be classified is collected in a training set defined here as,

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

The task is then to train a machine to learn the relationship between the data and their respective labels. This amounts to learning the geometrical structure of the space in which the two classes of data lie. In Support Vector Machines, the idea is to define a boundary separating these two classes as much as possible. This can be interpreted as 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. Maximizing the margin of the hyperplane is then equivalent to maximizing the distance between the class boundaries. Vapnik suggests that the form of the hyperplane be chosen from family of functions with sufficient capacity[1]. In particular,  $F$  contains functions for the linearly and non-linearly separable hyperplanes;

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

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

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  $\phi: \mathbf{x} \subset \mathcal{R}^n \rightarrow \mathcal{R}^m$  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;

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

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

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

The size of the soft-margin is governed by the positive slack variables,  $\xi_i$  which relax the conditions in (4). 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 SVM margin maximization problem is then formulated as;

$$\text{P.1} \left\{ \begin{array}{l} \text{minimize } \mathfrak{J}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^l w_i^2 + C \sum_{i=1}^l \xi_i \\ \text{subject to } \begin{cases} y_i (\sum_{i=1}^l w_i \varphi(\mathbf{x}_i) + b) \geq 1 - \xi_i \\ \xi_i > 0, \forall i = 1..l \end{cases} \end{array} \right. \quad (6)$$

The parameter C can be interpreted as a regularization parameter which controls the tradeoff between generalization accuracy and the number of training errors. A larger C generally results in lower training errors but poorer prediction accuracy.

The program in P.1 is difficult to solve using practical optimization methods. However, duality theory allows us to transform a program to an easier and more practical version to solve. Duality in optimization theory was mainly developed by Fenchel, Kuhn and Tucker, Rockafeller, Dorn, Wolfe and many others[8-11]. Their theorems mainly contributed to the understanding of the popular and important Lagrange Theory which describes how to obtain an equivalent dual program of a mathematical program. In summary, suppose we have a mathematical program of the following form;

$$\text{A.I} \left\{ \begin{array}{l} \text{minimize } \mathfrak{J}(\mathbf{w}) \\ \text{subject to constraints } \begin{cases} H(\mathbf{w}) \leq 0 \\ G(\mathbf{w}) = 0 \end{cases} \\ \text{where } H(\mathbf{w}) = \begin{pmatrix} h_1(\mathbf{w}) \\ \vdots \\ h_m(\mathbf{w}) \end{pmatrix} \quad G(\mathbf{w}) = \begin{pmatrix} g_1(\mathbf{w}) \\ \vdots \\ g_n(\mathbf{w}) \end{pmatrix} \end{array} \right. \quad (7)$$

If the objective function  $\mathfrak{J}(\mathbf{w})$ , equality constraints  $g_i(\mathbf{w})$  and inequality constraints  $h_i(\mathbf{w})$  all have smooth first partial derivatives which are linearly independent, there exists scalars  $\alpha, \beta \geq 0$  known as Lagrange Multipliers so that we can write;

$$\text{A.II} \left\{ \begin{array}{l} \text{minimize } L(\mathbf{w}) = f(\mathbf{w}) + \sum_{i=1}^n \beta_i g_i(\mathbf{w}) + \sum_{j=1}^m \alpha_j h_j(\mathbf{w}) \\ \text{subject to } \begin{cases} \nabla f(\mathbf{w}) + \sum_{i=1}^n \beta_i \nabla g_i(\mathbf{w}) + \sum_{j=1}^m \alpha_j \nabla h_j(\mathbf{w}) = 0 \\ \alpha, \beta > 0 \end{cases} \end{array} \right. \quad (8)$$

The program A.II is the *dual* to the program A.I and the solution to it is also a solution to A.I. Kuhn and Tucker later generalized this duality to include non-differentiable functions by introducing a saddle point theorem. It turns out that when the objective functions are non-differentiable, any saddle point of A.II is also a solution to A.I. We now apply Lagrange Theory to solve (6) giving us the Lagrange Primal problem of the following form;

$$\text{P.2} \left\{ \begin{array}{l} \text{minimize } \mathfrak{J}(w, \alpha) = \frac{1}{2} \sum_{i=1}^l w_i^2 + C \sum_{i=1}^l \xi_i + \sum_{i=1}^l \alpha_i (y_i (\sum_{j=1}^l w_j \varphi(x_j) + b) - 1 + \xi_i) \\ \quad + \sum_{i=1}^l \eta_i \xi_i \\ \text{subject to } \left\{ \begin{array}{l} \nabla \left( \frac{1}{2} \sum_{i=1}^l w_i^2 + C \sum_{i=1}^l \xi_i \right) + \nabla \left( \sum_{i=1}^l \alpha_i (y_i (\sum_{j=1}^l w_j \varphi(x_j) + b) - 1 + \xi_i) \right) \\ \quad + \sum_{i=1}^l \eta_i \xi_i \end{array} \right\} = 0 \\ \alpha_i \geq 0, \forall i=1..l \end{array} \right. \quad (9)$$

The program P.2 is popularly solved through the use of its dual representation which is found by incorporating the gradient condition in P.2 into the objective function and eliminating the primal variables. The stationary partial gradients of P.2 are given by;

$$\begin{aligned}
 \nabla_{\mathbf{w}} \mathfrak{J}(w, \alpha, \xi, b) &= \mathbf{w} + \sum_{i=1}^l \alpha_i y_i \varphi(\mathbf{x}_i) = \mathbf{0} \\
 \nabla_{\alpha} \mathfrak{J}(w, \alpha, \xi, b) &= \mathbf{y} \left( \sum_{i=1}^l w_i \varphi_i(\mathbf{x}) + b \right) - 1 + \boldsymbol{\xi} = 0 \\
 \nabla_b \mathfrak{J}(w, \alpha, \xi, b) &= \sum_{i=1}^l \alpha_i y_i = 0 \\
 \nabla_{\boldsymbol{\xi}} \mathfrak{J}(w, \alpha, \xi, b) &= C \boldsymbol{\xi} - \boldsymbol{\alpha} = \mathbf{0}
 \end{aligned} \quad (10)$$

We then eliminate the primal variables through back substitution, giving us the Lagrange Dual;

$$\text{P.3} \left\{ \begin{array}{l} \text{minimize } \mathfrak{J}(\boldsymbol{\alpha}) = -\sum_{i=1}^l \alpha_i + \frac{1}{2} \sum_{i,j=1}^l \alpha_i \alpha_j y_i y_j \langle \varphi(\mathbf{x}_i), \varphi(\mathbf{x}_j) \rangle \\ \text{subject to } \left\{ \begin{array}{l} 0 \leq \alpha_i \leq C \\ \sum_{i=1}^l \alpha_i y_i = 0 \end{array} \right. \end{array} \right. \quad (11)$$

In P.3 we find that the feature vectors exist as dot products and can be represented by a kernel function using Mercer's Theorem. The formulation avoids the task of explicitly specifying the feature space mapping and we simply have to choose a kernel function;

$$K(\mathbf{x}_i, \mathbf{x}_j) = \langle \varphi(\mathbf{x}_i), \varphi(\mathbf{x}_j) \rangle \quad (12)$$

When solving the dual programs P.3, we can derive the following dual hyperplane form expressed solely in terms of the Lagrange Multipliers by using the gradient condition in (10) to substitute for (4);

$$f(x) = \sum_{i=1}^l \alpha_i y_i K(x_i, x) + b \quad (13)$$

The SVM solution is the vectors  $w^*, \alpha^*$  which satisfy P.1-P.3 giving the trained Support Vector Classifier defined by the decision functions;

$$f(x) = \text{sign} \left( \sum_{i=1}^l w_i^* x_i + b \right) \quad (14)$$

$$f(x) = \text{sign}\left(\sum_{i=1}^l \alpha_i^* y_i K(x_i, x) + b\right) \quad (15)$$

The asterisk denotes the optimal values of the variables and we should state that the form of (14) is only practical for linearly separable Support Vector Machine classifiers. It should be noted that the solution is sparse, meaning that a lot of  $\alpha_i = 0$  and the decision function,  $f(x)$  could be represented solely by the Support Vectors, (i.e.  $\alpha_i \neq 0$ ).

### III. A PRACTICAL LINEAR PROGRAM: MINIMIZATION OF TRAINING ERRORS

The dual form of P.3 is in essence a constrained quadratic program which requires optimizers that use activesets, interior point methods and so on. However, simpler programs have been devised to solve a different problem to P.3 choosing instead to solve the following program;

$$\text{P.4} \begin{cases} \text{minimize } \nabla \mathfrak{Z}(\boldsymbol{\alpha})^T d \\ \text{where } \begin{cases} 0 \leq \alpha_i + d_i \leq C \\ \sum_{i=1}^l y_i d_i = 0 \end{cases} \end{cases} \quad (16)$$

Their algorithms terminate using a reformulated Karush-Kuhn-Tucker (KKT) conditions[12];

$$\begin{aligned} \nabla \mathfrak{Z}(\boldsymbol{\alpha}^*)_i + by_i &\geq 0 & \alpha_i &= 0 \\ \nabla \mathfrak{Z}(\boldsymbol{\alpha}^*)_i + by_i &\leq 0 & \alpha_i &= C \\ \nabla \mathfrak{Z}(\boldsymbol{\alpha}^*)_i + by_i &= 0 & 0 < \alpha_i < C \end{aligned} \quad (17)$$

Most of the successful SVM optimization algorithms mentioned before, including the popular SMO and SVMlight further use a decomposition method, first proposed by Osuna[13] to tackle the problem of memory with large datasets. However [12] have mentioned that decomposition methods tend to be slower than Newtonian methods. Unfortunately, Newtonian methods require that the entire kernel matrix be kept in memory which may not be possible for datasets that have 10000 points or more. Thus a compromise on which algorithm to use when solving a SVM problem has to be made. Therefore, the next logical thing which follows is to investigate a possible design of a hybrid optimization algorithm that possibly uses some Newtonian or linear iterative update method combined with a form of space decomposition to solve P.4. The problem with P.4 is that it is not immediately clear how to apply linear iterative methods which are specifically designed to solve a set of linear equations; in particular partial differential equations. The motive for this paper is then clear, for to apply these iterative methods we have to first establish a system of linear equations closely related to P.4. In the following, we show how this is done by deriving (16) from basic principles first and then forming the set of linear equations in matrix form.

#### A. THE GRADIENT OF THE UNCONSTRAINED QUADRATIC PROGRAM

In this section, we first investigate the **unconstrained** QP problem of P.3 derived by dropping the constraints on the objective function. We can easily rewrite the function  $\psi(\boldsymbol{\alpha})$  in matrix form which we now define as;

$$\begin{aligned}\psi(\boldsymbol{\alpha}) &= \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i,j=1}^l \alpha_i \alpha_j y_i y_j K(x_i, x_j) \\ &= \boldsymbol{\alpha}^T \mathbf{1} - \frac{1}{2} \boldsymbol{\alpha}^T \mathbf{G} \boldsymbol{\alpha}\end{aligned}\quad (18)$$

$$\text{where } \begin{cases} \boldsymbol{\alpha}, \mathbf{1} \in \mathfrak{R}^l \\ \mathbf{G} \in \mathfrak{R}^l \times \mathfrak{R}^l \\ G_{ij} = y_i y_j K(x_i, x_j) \\ y_i \in (1, -1) \quad \forall i = 1..l \end{cases}$$

The function  $\psi(\boldsymbol{\alpha})$  is quadratic in terms of the variable  $\alpha$  (Lagrange Multiplier) and has a unique maximum which we denote as  $\boldsymbol{\alpha}^*$ . As before the asterisk defines the optimal values of the  $\alpha$  vector, which in this case refers to the unique maximizer. The maximum occurs when all partial derivatives vanish or simply the gradients with respect to the vector  $\boldsymbol{\alpha}$  become stationary, i.e

$$\nabla \psi(\boldsymbol{\alpha}^*) = 0 \quad (19)$$

For a quadratic function, this condition is necessary and sufficient (ref kaplan) in order to obtain the optimal vector  $\boldsymbol{\alpha}^*$  which is the unique maximizer to  $\psi(\alpha)$ . The gradients of (18) with respect to the Lagrange Multipliers,  $\alpha_i$  can be found by taking the partial derivatives of (18) which is Gateux-differentiable giving;

$$\nabla_{\alpha_i} \psi(\alpha) = 1 - y_i \sum_{j=1}^l \alpha_j y_j K(x_i, x_j) \quad \forall i = 1..l \quad (20)$$

We verify the differentiability in the following simple example for  $\boldsymbol{\alpha} \in \mathfrak{R}^2$ . Let us denote  $K_{ij} = K(x_i, x_j)$  and expanding the function  $\psi(\boldsymbol{\alpha})$ , we get;

$$\begin{aligned}\psi(\boldsymbol{\alpha}) &= \alpha_1 + \alpha_2 - \frac{1}{2} (\alpha_1^2 y_1^2 K_{11} + \alpha_1 \alpha_2 y_1 y_2 K_{12} + \alpha_2 \alpha_1 y_2 y_1 K_{21} + \alpha_2^2 y_2^2 K_{22}) \\ &= \alpha_1 + \alpha_2 - \frac{1}{2} (\alpha_1^2 y_1^2 K_{11} + 2\alpha_1 \alpha_2 y_1 y_2 K_{12} + \alpha_2^2 y_2^2 K_{22}) \\ \nabla_{\alpha_1} \psi(\boldsymbol{\alpha}) &= \frac{\partial \psi(\boldsymbol{\alpha})}{\partial \alpha_1} = 1 - \frac{1}{2} (2\alpha_1 y_1^2 K_{11} + 2\alpha_2 y_1 y_2 K_{12}) \\ &= 1 - y_1 (\alpha_1 y_1 K_{11} + \alpha_2 y_2 K_{12}) \\ &= 1 - y_1 \sum_{j=1}^2 \alpha_j y_j K_{1j}\end{aligned}\quad (21)$$

We can clearly obtain (21) from (20) if we set  $i = 1$  and we can quickly check that this is true also for  $i = 2$ . We can now generalize this to the entire vector  $\boldsymbol{\alpha}$  and compute the maximizer  $\boldsymbol{\alpha}^*$  directly by applying (19);

$$\begin{aligned}\psi(\boldsymbol{\alpha}) &= \boldsymbol{\alpha}^T \mathbf{1} - \frac{1}{2} \boldsymbol{\alpha}^T \mathbf{G} \boldsymbol{\alpha} \\ \nabla_{\boldsymbol{\alpha}} \psi(\boldsymbol{\alpha}) &= \mathbf{1}^T \mathbf{1} - \frac{1}{2} \mathbf{1}^T \mathbf{G} \boldsymbol{\alpha} = 0\end{aligned}\quad (22)$$

The unique maximizer is then found by;

$$\begin{aligned}\rightarrow \mathbf{1}^T \mathbf{2} &= \mathbf{1}^T \mathbf{G} \boldsymbol{\alpha} \\ \therefore \boldsymbol{\alpha}^* &= \mathbf{G}^{-1} \mathbf{2}\end{aligned}$$

Resubstituting this back into (18), provided  $\mathbf{G}$  is non-singular, we get the maximum value of  $\psi(\mathbf{a})$  as;

$$\begin{aligned}\psi(\mathbf{a}^*) &= (\mathbf{G}^{-1}\mathbf{2})^T \mathbf{1} - \frac{1}{2}(\mathbf{G}^{-1}\mathbf{2})^T \mathbf{G}\mathbf{G}^{-1}\mathbf{2} \\ &= (\mathbf{G}^{-1}\mathbf{2})^T \mathbf{1} - (\mathbf{G}^{-1}\mathbf{2})^T \mathbf{1} = 0\end{aligned}\quad (23)$$

This is somewhat interesting because we can conclude that no matter what the size of the training data set, the maximum value possible for P.3 is zero provided the matrix  $\mathbf{G}$  is nonsingular. The constraints in P.3 actually restrict the feasible region of solutions which results in P.3 having a value larger than zero. A non-singular  $\mathbf{G}$  is guaranteed by using a positive definite kernel, e.g gaussian kernel and provided there are no duplicate examples with opposite labels. We will elaborate on the possibilities further arising from this observation in the next section.

It can be shown quickly that if minimize the negative of  $\psi(\mathbf{a})$ , we will end up with the same unique solution provided the condition of  $\det \mathbf{G} \neq 0$  on the matrix  $\mathbf{G}$  holds. This fact is well known and, we have a trivial dual of the form;

$$\psi(\mathbf{a}) \underset{\max}{=} -\psi(\mathbf{a}) \underset{\min}{=} \quad (24)$$

In fact the negative maximization of P.3 is known as the Wolfe dual in the literature (ref smola). We now proceed further with a minor observation. Scalar multiples of the function  $\psi(\mathbf{a})$  gives the same unique maximizer too. We show this in the following technical lemma.

**Lemma III.1:** Let the quadratic function  $\psi(\mathbf{a})$  where  $\mathbf{a} \in \mathfrak{R}^l$ , have a unique maximum denoted by  $\mathbf{a}^*$ . Then for any nonzero scalar  $k, k \in \mathfrak{R}$  the function  $k\psi(\mathbf{a})$  possesses the same unique maximizer solution  $\mathbf{a}^*$ .

**Proof:**

$$\text{Let } \psi(\mathbf{a}) = \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i,j=1}^l \alpha_i \alpha_j y_i y_j K(x_i, x_j) \quad (25)$$

Then

$$k\psi(\mathbf{a}) = k \sum_{i=1}^l \alpha_i - \frac{k}{2} \sum_{i,j=1}^l \alpha_i \alpha_j y_i y_j K(x_i, x_j) \quad (26)$$

The solution to (26) is found as before when all gradients with respect to the variables are stationary. In particular, we require that  $\forall i = 1..l$ ;

$$\begin{aligned}\nabla_{\alpha_i} k\psi(\mathbf{a}) &= \nabla \left( k \sum_{i=1}^l \alpha_i - \frac{k}{2} \sum_{i,j=1}^l \alpha_i \alpha_j y_i y_j K(x_i, x_j) \right) \\ &= k - \frac{ky_i}{2} \sum_{i,j=1}^l \alpha_j y_j K(x_i, x_j) \\ &= k \left( 1 - \frac{y_i}{2} \sum_{i,j=1}^l \alpha_j y_j K(x_i, x_j) \right) \\ &= k \nabla_{\alpha_i} \psi(\mathbf{a})\end{aligned}$$

At the maximum point,  $k \nabla \psi(\mathbf{a}) = 0$  which holds if and only if  $\nabla \psi(\mathbf{a}) = 0, \forall \alpha_i \in \mathbf{a}$  and hence the solution is also  $\mathbf{a}^*$   $\square$

The general element gradient function, (20) by itself seems rather uninteresting, yet it is explicitly used as the objective function for the program (16). We know that they have to be zero at the maximum



or minimum points of the quadratic function and this is what optimization methods such as gradient descent methods iterate on. However, keeping in mind that we are examining the objective function of P.3 alone minus the constraints, our previous observation gives us the idea that we could manipulate scalar multiples of the gradients of a quadratic function to a form which we could associate easily with. We now use the result in **Lemma III.1** to tie in with a well-known quantity in the SVM classification problem by constructing the following theorem.

**Theorem III.1:** Let  $\psi(\mathbf{a})$  be an unconstrained SVM quadratic function having the form,

$$\psi(\mathbf{a}) = \mathbf{a}^T \mathbf{1} - \frac{1}{2} \mathbf{a}^T \mathbf{G} \mathbf{a}$$

where the elements of  $\mathbf{G}$  are defined as in (18). Let the unique maximum be denoted by  $\mathbf{a}^*$ . Then for  $\mathbf{a} \neq \mathbf{a}^*$  there exists a scalar,  $\delta_i$  i.e  $\delta_i \geq 0$  such that;

$$y_i - f_E(x_i, \mathbf{a}) \leq |\delta_i| \quad \forall i = 1..l$$

where  $f_E(x_i, \mathbf{a}) = \sum_{j=1}^l \alpha_j y_j K(x_i, x_j) + b$  and  $\lim_{\delta_i \rightarrow 0} \mathbf{a} = \mathbf{a}^*, \forall i = 1..l$  holds whenever  $\mathbf{G}$  is nonsingular.

**Proof:**

Let

$$v(\alpha) = y_i - \sum_{j=1}^l \alpha_j y_j K(x_i, x_j) \quad (27)$$

By construction, the vector  $\mathbf{a}$  lies in the convex set  $\mathfrak{R}^l$ . Then by the properties of a convex set, for any maximal  $\mathbf{a}^*$ , there exists a vector,  $\alpha$  in  $\mathfrak{R}^l$  such that the vector  $z = (1-t)\mathbf{a}^* + t\alpha$  also lies in  $\mathfrak{R}^l$  for some  $0 < t \leq 1$ . We then have;

$$\begin{aligned} v(z) &= v\left((1-t)\mathbf{a}^* + t\alpha\right) \\ &= y_i - \sum_{i,j=1}^l \left((1-t)\alpha_j^* + t\alpha_j\right) y_j K(x_i, x_j) \\ &= y_i - \sum_{i,j=1}^l \alpha_j^* y_j K(x_i, x_j) + t \sum_{i,j=1}^l \alpha_j^* y_j K(x_i, x_j) - t \sum_{i,j=1}^l \alpha_j y_j K(x_i, x_j) \\ &= v(\alpha^*) + t \sum_{i,j=1}^l (\alpha_j^* - \alpha_j) y_j K(x_i, x_j) \\ &\leq v(\alpha^*) + t \sum_{i,j=1}^l \alpha_j^* y_j K(x_i, x_j) \end{aligned} \quad (28)$$

Let  $t = 1$ , we can see that the right hand most term is a constant which we can denote by  $\varepsilon$  and we now rewrite the inequality as ;

$$v(\alpha) - v(\alpha^*) \leq |\varepsilon| \quad (29)$$

Now, we can see from (27) and (20) that by construction;

$$v(\alpha) = y_i \nabla_{\alpha_i} \psi(\alpha)$$

Then by applying **Lemma III.1**, we then have  $v(\alpha^*) = 0$  and using (13) and (19), we have;

$$\begin{aligned} v(\alpha) - v(\alpha^*) &= y_i - \sum_{j=1}^l \alpha_j y_j K(x_i, x_j) \\ &= y_i - (f_E(x_i, \mathbf{a}) - b) \\ &= y_i - f_E(x_i, \mathbf{a}) + b \leq |\varepsilon_i| \end{aligned}$$

$$\rightarrow y_i - f_E(x_i, \mathbf{a}) \leq |\varepsilon_i - b| \quad (30)$$

We can then set  $|\delta_i| = |\varepsilon_i - b|$  and see from (28) that  $\delta_i \rightarrow 0, \forall i = 1..l$  implies that  $\alpha \rightarrow \alpha^*$  and thus the limit is proved. If  $\mathbf{G}$  is singular, then  $\psi(\mathbf{a})$  has a critical point which is non-unique. Hence the theorem holds only if  $\mathbf{G}$  is nonsingular.  $\square$

The quantity  $y_i - f_E(x_i, \mathbf{a})$  is better known as the error on a particular training example, or simply the training error usually denoted by  $E_i$ . We refer to this as the dual form of the training error, because incidentally if we examine our discriminant function (5) from the original problem closely, we recover the following;

$$\begin{aligned} & y_i \left( \sum_{i=1}^n w_i \varphi(\mathbf{x}_i) + b \right) \geq 1 \\ \rightarrow & y_i \left( \sum_{i=1}^n w_i \varphi(\mathbf{x}_i) + b \right) - 1 = \left| \sum_{i=1}^n w_i \varphi(\mathbf{x}_i) + b - y_i \right| \geq 0 = \left| y_i - \sum_{i=1}^n w_i \varphi(\mathbf{x}_i) + b \right| \leq 0 \end{aligned}$$

This quantity happens to be the training error expressed in terms of the original hyperplane, which we believe could be referred to fondly as the primal training error form. The results in the **Theorem III.1** show that if we set the dual training error for each example to zero, we would obtain the unique maximizer,  $\mathbf{a}^*$  to the quadratic function  $\psi(\mathbf{a})$ . In fact, ensuring zero training error using the dual hyperplane form (13) is now shown to be implicitly equivalent to ensuring that the partial gradients of  $\psi(\alpha)$  are stationary. Recall that the dual hyperplane form also ensures indirectly that the margin's are maximized in the weight space,  $\mathbf{w}$ .

### B. The Linear Program

We can now construct a constrained linear program by applying the same set of constraints and minimizing the sum of the dual training errors instead. We postulate that this is different from direct Empirical Risk Minimization because our use of the dual hyperplane form implicitly enforces capacity control through minimization of the hyperplane margins. We now have the following mathematical program;

$$\text{P.5} \begin{cases} \text{minimize} & \mathfrak{Z}(\mathbf{a}) = \sum_{i=1}^l \left| y_i - \left( \sum_{j=1}^l \alpha_j y_j K(x_i, x_j) + b \right) \right| \\ \text{subject to} & \begin{cases} 0 \leq \alpha_i \leq C \\ \sum_{i=1}^l \alpha_i y_i = 0 \end{cases} \end{cases} \quad (31)$$

In fact, (30) in **Theorem III.1** allows us to reconcile (31) with the previous linear program described in Joachims etc.

Rewriting (31), we obtain

$$\begin{aligned} \mathfrak{Z}(\alpha) &= \sum_{i=1}^l \left| y_i - \left( \sum_{j=1}^l \alpha_j y_j K(x_i, x_j) + b \right) \right| \\ &= \sum_{i=1}^l y_i \left| 1 - \left( \sum_{j=1}^l \alpha_j y_j K(x_i, x_j) \right) \right| + b \end{aligned}$$

Each element then satisfies;

$$\begin{aligned}
 y_i \mathfrak{S}(\alpha) &\rightarrow \sum_{i=1}^l y_i \nabla \psi_i(\mathbf{a}) - y_i b \leq y_i | \varepsilon_i - b | \\
 &\rightarrow \sum_{i=1}^l y_i \nabla \psi_i(\mathbf{a}) \leq y_i | \varepsilon_i | \\
 &\rightarrow \sum_{i=1}^l y_i \nabla \psi_i(\mathbf{a}) - y_i | \varepsilon_i | \leq 0 \\
 &\rightarrow \sum_{i=1}^l d_i \nabla \psi_i(\mathbf{a}) \leq 0 \\
 &\rightarrow \nabla \psi(\mathbf{a})^T d \leq 0
 \end{aligned}$$

since  $\varepsilon_i$  is the constant from **Theorem III.1** and equality is achieved at optimal  $\alpha$ . We can now write the problem as a system of linear equations solved on a compact domain;

$$\begin{aligned}
 Hu &= F \\
 \text{subject to: } u &\in Q
 \end{aligned} \tag{32}$$

$$H = \begin{pmatrix} K_{11} & \cdots & K_{1l} \\ \vdots & \ddots & \vdots \\ K_{l1} & \cdots & K_{ll} \end{pmatrix} \quad \mathbf{u} = \begin{pmatrix} y_1 \alpha_1 \\ \vdots \\ y_l \alpha_l \end{pmatrix} \quad \mathbf{F} = \begin{pmatrix} y_1 - b \\ \vdots \\ y_l - b \end{pmatrix}$$

The matrix H is better known as the Gram matrix, which is symmetrical, and positive definite provided we use a positive definite kernel. The vector of variables  $\mathbf{u}$  could be written solely, in terms of the vector  $\alpha$  but this would result in a non-symmetrical matrix H which would be unsuitable for acceleration methods suggested in [14]. The region Q is the feasible region of solutions, which will be discussed further in the next section. It is appropriate at this point to note that the iterates are bounded by this region, and hence optimization speeds using gradient methods or step length searches should have different rates of convergence.

We can further incorporate the equality constraint into the objective function through methods of elimination. This will be useful for writing the problem in matrix form without the step of computing the threshold, b. We note, that the problem formulation commonly used in [5, 15] treats the threshold as a Lagrange Multiplier which allows the problem to be formed as a slightly different set of linear equations. We first apply some elimination methods in the following sections.

#### 1) Direct Elimination

This is intuitively the simplest method which stems from solving equations with multiple variables in linear algebra. Consider the equality constraint expanded as a single equation. We now have the following;

$$\alpha_1 y_1 + \alpha_2 y_2 + \dots + \alpha_l y_l = 0 \tag{33}$$

where the elements  $\alpha_j \in \alpha, y_j \in \mathbf{y}, \forall j = 1 \dots l$ . Let us rewrite each element in terms of the others, we obtain;

$$\begin{aligned}
 \alpha_1 &= \frac{1}{y_1} (\alpha_2 y_2 + \alpha_3 y_3 + \dots + \alpha_l y_l) \\
 \alpha_2 &= \frac{1}{y_2} (\alpha_1 y_1 + \alpha_3 y_3 + \dots + \alpha_l y_l) \\
 &\vdots \\
 \alpha_l &= \frac{1}{y_l} (\alpha_1 y_1 + \alpha_2 y_2 + \dots + \alpha_{l-1} y_{l-1})
 \end{aligned}$$

Then a particular Lagrange Multiplier can be expressed in the following relationship;

$$\alpha_j = \frac{1}{y_j} \sum_{\substack{k=1 \\ k \neq j}}^l \alpha_k y_k \quad (34)$$

We can now substitute (34) back into the equation of the objective function P.5 to give;

$$\begin{aligned} \mathfrak{J}(\boldsymbol{\alpha}) &= \sum_{i=1}^l \left[ y_i - \left( \sum_{j=1}^l \alpha_j y_j K(x_i, x_j) + b \right) \right] \\ &\rightarrow \sum_{i=1}^l \left[ y_i - \left( \sum_{j=1}^l \frac{1}{y_j} \sum_{\substack{k=1 \\ k \neq j}}^l \alpha_k y_k y_j K(x_i, x_j) + b \right) \right] \\ &\rightarrow \sum_{i=1}^l \left[ y_i - \left( \sum_{\substack{j=1 \\ k \neq j}}^l \sum_{k=1}^l \alpha_k y_k K(x_i, x_j) + b \right) \right] \\ &\rightarrow \sum_{i=1}^l \left[ y_i - \left( \sum_{j=1}^l \left( \sum_{k=1}^l \alpha_k y_k - \alpha_j y_j \right) K(x_i, x_j) + b \right) \right] \end{aligned}$$

We now obtain the following linear program.

$$\text{P.6} \begin{cases} \text{minimize } \mathfrak{J}(\boldsymbol{\alpha}) = \sum_{i=1}^l \left[ y_i - \left( \sum_{j=1}^l \left( \sum_{k=1}^l \alpha_k y_k - \alpha_j y_j \right) K(x_i, x_j) + b \right) \right] \\ \text{subject to } \{ 0 \leq \alpha_i \leq C \end{cases} \quad (35)$$

## 2) Generalized Elimination Method

The direct elimination method might not be the best method or even practical to implement. A generalized elimination method[8] exists which uses a linear transformation on the variable, in this case the Lagrange Multiplier to reduce a problem with equality constraints to an unconstrained problem. In our case, we show that this method reduces to the direct elimination method. We first note that the equality constraint can be written in vector form as;

$$\mathbf{y}^T \boldsymbol{\alpha} = 0 \quad (36)$$

Let  $\mathbf{A}$  and  $\mathbf{B}$  be a  $l \times 1$  vector and a  $l \times (l-1)$  matrix respectively such that  $[A : B]$  is a non-singular matrix ;

$$[A : B] = \begin{bmatrix} a_1 & b_{11} & \cdots & b_{1,l-1} \\ \vdots & \vdots & \vdots & \vdots \\ a_l & b_{l1} & \cdots & b_{l,l-1} \end{bmatrix}$$

Furthermore, let us choose  $\mathbf{A}$  and  $\mathbf{B}$  so that  $\mathbf{Y}^T \mathbf{A} = 1$  and  $\mathbf{Y}^T \mathbf{B} = \mathbf{0}$ .  $\mathbf{A}$  is generally a  $l \times k$  matrix if there happens to be  $k$ -equality constraints, but in our case it is a vector due to the single equality constraint. In any case, let us treat  $\mathbf{A}$  as a generalized left inverse for  $\mathbf{Y}$  so that we can say that;

$$\begin{aligned} \mathbf{Y}^T \boldsymbol{\alpha} &= 0 \\ \rightarrow \boldsymbol{\alpha} &= \mathbf{A}^T \mathbf{0} = \mathbf{0} \end{aligned}$$

We state that this is simply the trivial solution of  $\boldsymbol{\alpha}$  and is of course not unique because we can define other feasible points by taking any feasible direction from the trivial solution, i.e

$$\boldsymbol{\alpha} = \mathbf{A}^T \mathbf{0} + \boldsymbol{\kappa} \quad (37)$$

provided of course that  $\mathbf{\kappa}$  is a direction in the linear space  $\{\mathbf{\kappa} | \mathbf{Y}^T \mathbf{\kappa} = 0\}$  with dimension  $l-1$ . Now, let the matrix  $\mathbf{B}$ , contain linearly independent columns,  $\{\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_{l-1}\}$  which can be viewed as reduced coordinate directions and the scalars  $\lambda_i$  form the reduced variables of these coordinate directions. Then any feasible direction,  $\mathbf{\kappa}$  is given by;

$$\mathbf{\kappa} = \mathbf{B}\boldsymbol{\lambda} = \sum_{i=1}^{l-1} \mathbf{b}_i \lambda_i \quad (38)$$

Fig III.1 shows an example of the feasible direction in terms of reduced coordinates when  $\boldsymbol{\alpha} \in \mathfrak{R}^3$ .

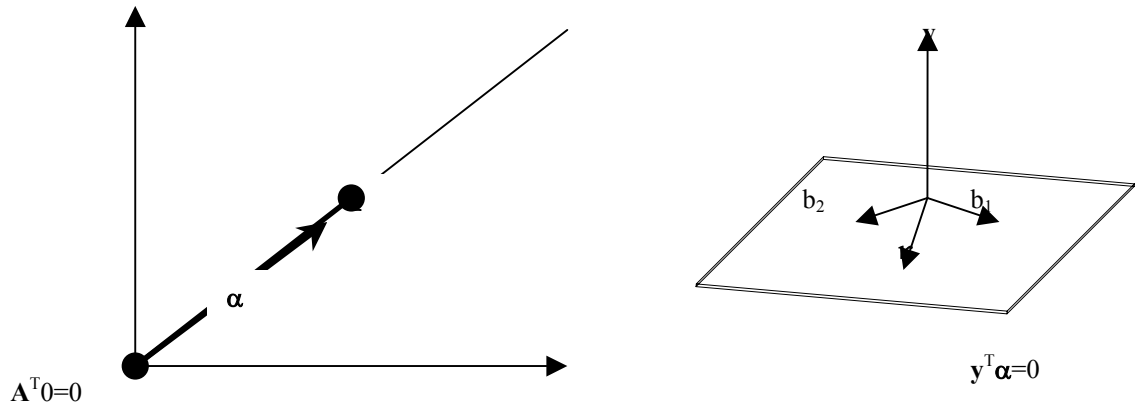


Fig III.1: The feasible points of  $\boldsymbol{\alpha}$  can be found by taking a step from the trivial point in the feasible direction  $\mathbf{\kappa}$ . The reduced coordinate vectors  $\mathbf{b}_1$  and  $\mathbf{b}_2$  determine the feasible direction,  $\mathbf{\kappa}$ .

The problem that remains now is to choose a suitable matrix  $\mathbf{B}$  which has linearly independent columns. We propose one such matrix  $\mathbf{B}$  derived from the condition  $\mathbf{Y}^T \mathbf{B} = \mathbf{0}$  having the following form;

$$\mathbf{B} = \begin{bmatrix} \alpha_2 & \alpha_3 & \dots & \alpha_l \\ \alpha_3 & \alpha_4 & \dots & \alpha_1 \\ \vdots & \vdots & \dots & \vdots \\ \alpha_1 & \alpha_2 & \dots & \alpha_{l-1} \end{bmatrix} \quad (39)$$

Simplifying (37), we now have the following  $\boldsymbol{\alpha}$  vector expressed in reduced coordinate directions as;

$$\boldsymbol{\alpha} = \mathbf{B}\boldsymbol{\lambda} \quad (40)$$

We now can substitute this relationship into (31) to give us the following program;

$$\text{P.7} \left\{ \begin{array}{l} \text{minimize } \mathfrak{S}(\boldsymbol{\alpha}) = \sum_{i=1}^l \left[ y_i - \left( \sum_{j=1}^l \left( \sum_{k=1, k \neq j}^{l-1} \lambda_k \alpha_k \right) y_j K(x_i, x_j) + b \right) \right] \\ \text{subject to } \left\{ \begin{array}{l} 0 \leq \alpha_k \leq \frac{C(l-1)}{\lambda_k} \end{array} \right. \end{array} \right. \quad (41)$$

We note that  $\mathbf{\kappa}$  is sometimes known as a feasible correction.

### 3) Geometrical Elimination

Our equality constraint is unique because we can also use simple plane geometry to incorporate it into the objective function. In our previous discussion, we treated the equality constraint as a hyperplane which had a section lying in a convex polyhedra. We now expand this hyperplane idea by rewriting (36) in the general hyperplane equation form;

$$\mathbf{y} \cdot (\boldsymbol{\alpha} - \mathbf{p}) = 0 \quad (42)$$

The normal to the hyperplane is the vector,  $\mathbf{y}$  and any point on the hyperplane satisfies equation (42). The point,  $\mathbf{p}$  is any point lying on the hyperplane surface. Since the hyperplane passes through the origin, then  $\mathbf{p} = \mathbf{0}$  is one such point. All we have to do now is to find the general parametric equation of the feasible alpha points lying on the plane and substitute it back into the objective function. We would have then incorporated the linear equality constraint into the objective function.

First, let us recall that a hyperplane in the space  $\mathfrak{R}^l$  is a  $l-1$ -dimensional affine subspace of  $\mathfrak{R}^l$ . We expand (42) to give ,

$$y_1(\alpha_1 - p_1) + y_2(\alpha_2 - p_2) + \dots + y_l(\alpha_l - p_l) = 0 \quad (43)$$

Now, we can choose any nonzero scalar,  $y_i$  to solve for the corresponding point  $(\alpha_i - p_i)$ . Since this is classification and the scalars are all either 1 or -1, we could choose any point. For the sake of this discussion, let us solve for the point  $\alpha_i$ . We get;

$$\alpha_i = \frac{y_1}{y_i}(\alpha_1 - p_1) - \frac{y_2}{y_i}(\alpha_2 - p_2) + \dots - \frac{y_{l-1}}{y_i}(\alpha_{l-1} - p_{l-1})$$

If we set ,  $\{\alpha_1, \alpha_2, \dots, \alpha_{l-1}\} = \{\lambda_1, \lambda_2, \dots, \lambda_{l-1}\}$  where  $-\infty < \lambda_i < \infty, \forall i = 1..l$ , we get the following;

$$\begin{aligned} \boldsymbol{\alpha} &= \{\alpha_1, \alpha_2, \dots, \alpha_l\} \\ &= \left\{ \lambda_1, \lambda_2, \dots, \lambda_{l-1}, -\frac{y_1}{y_l}(\lambda_1 - p_1) - \frac{y_2}{y_l}(\lambda_2 - p_2) + \dots - \frac{y_{l-1}}{y_l}(\lambda_{l-1} - p_{l-1}) \right\} \\ &= \lambda_1 \begin{pmatrix} 1 \\ 0 \\ \vdots \\ -\frac{y_1}{y_l} \left( \lambda_1 - \frac{p_1}{\lambda_1} \right) \end{pmatrix} + \lambda_2 \begin{pmatrix} 0 \\ 1 \\ \vdots \\ -\frac{y_2}{y_l} \left( \lambda_2 - \frac{p_2}{\lambda_2} \right) \end{pmatrix} + \dots + \lambda_{l-1} \begin{pmatrix} 0 \\ \vdots \\ 1 \\ -\frac{y_{l-1}}{y_l} \left( \lambda_{l-1} - \frac{p_{l-1}}{\lambda_{l-1}} \right) \end{pmatrix} \end{aligned} \quad (44)$$

We can now simply substitute this relationship into (31) to get the objective function as;

$$\begin{aligned}
 \mathfrak{Z}(\boldsymbol{\alpha}) &= \sum_{i=1}^l \left[ y_i - \left( \sum_{j=1}^l \alpha_j y_j K(x_i, x_j) + b \right) \right] \\
 &\rightarrow \sum_{i=1}^l \left[ y_i - \left( \sum_{j=1}^{l-1} \lambda_j y_j K(x_i, x_j) - \sum_{k=1}^{k-1} \frac{y_k}{y_l} \left( \lambda_k - \frac{P_k}{\lambda_k} \right) y_l K(x_i, x_l) + b \right) \right] \\
 &\rightarrow \sum_{i=1}^l \left[ y_i - \left( \sum_{j=1}^{l-1} \lambda_j y_j K(x_i, x_j) - \sum_{k=1}^{k-1} y_k \left( \lambda_k - \frac{P_k}{\lambda_k} \right) K(x_i, x_l) + b \right) \right] \\
 &\rightarrow \sum_{i=1}^l \left[ y_i - \left( \sum_{j=1}^{l-1} \lambda_j y_j K(x_i, x_j) - \sum_{k=1}^{k-1} \lambda_k y_k K(x_i, x_l) + b \right) \right]
 \end{aligned}$$

We obtain the last step by using the origin as a known point,  $\mathbf{p}$  on the hyperplane. We then obtain the linear program as;

$$\text{P.8} \begin{cases} \text{minimize } \mathfrak{Z}(\boldsymbol{\alpha}) = \sum_{i=1}^l \left[ y_i - \left( \sum_{j=1}^{l-1} \lambda_j y_j K(x_i, x_j) - \sum_{k=1}^{k-1} \lambda_k y_k K(x_i, x_l) + b \right) \right] \\ \text{subject to } \{ 0 \leq \lambda_k \leq C \end{cases} \quad (45)$$

Let us make a note that this is closely similar in nature to the results of the direct elimination and the generalized elimination method. In fact, P6-P8 are the same linear program obtained through different methods of elimination. The geometrical method gives us a more general form of the linear program in terms of known points  $\mathbf{p}$  which are also feasible solutions  $\boldsymbol{\alpha}$ . This will be particularly useful for the design of a geometrical algorithm which solves the SVM problem based on search directions on the hyperplane enforced by the equality constraint.

#### IV. GEOMETRICAL VIEW OF SVM SOLUTIONS

##### A. Feasible Region of Solutions

The solution to the P.3-P.8, must lie in the feasible region,  $\mathbb{Q}$  defined by the intersections of the constraint set. We denote the region algebraically as follows;

$$\begin{aligned}
 \mathbb{Q} &= H(\boldsymbol{\alpha}) \cap \bigcap_{i=1}^l g(\alpha_i) \\
 g(\boldsymbol{\alpha}) &= \sum_{i=1}^l \alpha_i y_i, H = (\alpha_j | (\alpha_j - C) \leq 0) \\
 &\text{where } i, j, \alpha_j \geq 0, C > 0
 \end{aligned} \quad (46)$$

The set H is of size  $l$ , and forms a  $l$ -sided convex polyhedra lying entirely in the positive quadrant of the  $\mathfrak{R}^l$  space with each element  $h_j$  forming a boundary plane. Recall that  $l$  is the number of training examples and hence directly influences the dimension of the polyhedra. The region dictated by the equality constraint  $g$  is a linear manifold or hyperplane intersecting the polyhedra passing through the origin. Our optimal solution,  $\boldsymbol{\alpha}^*$  must lie in the region  $\mathbb{Q}$  or in other words the optimal vector of Lagrange Multipliers is a convex subset of  $\mathbb{Q}$ , i.e  $\boldsymbol{\alpha}^* \in \mathbb{Q}$ .

We illustrate this in the following example, for  $l=3$ , we have the following 3-dimension polyhedra having the form of a cube in Fig IV.1 below. The linear manifold is now the plane

$g(\alpha) = \sum_{i=1}^3 \alpha_i y_i$  which intersects the cube and the feasible region,  $\mathbb{Q}$  of solutions is now restricted to the surface of the plane contained within the region of the cube.

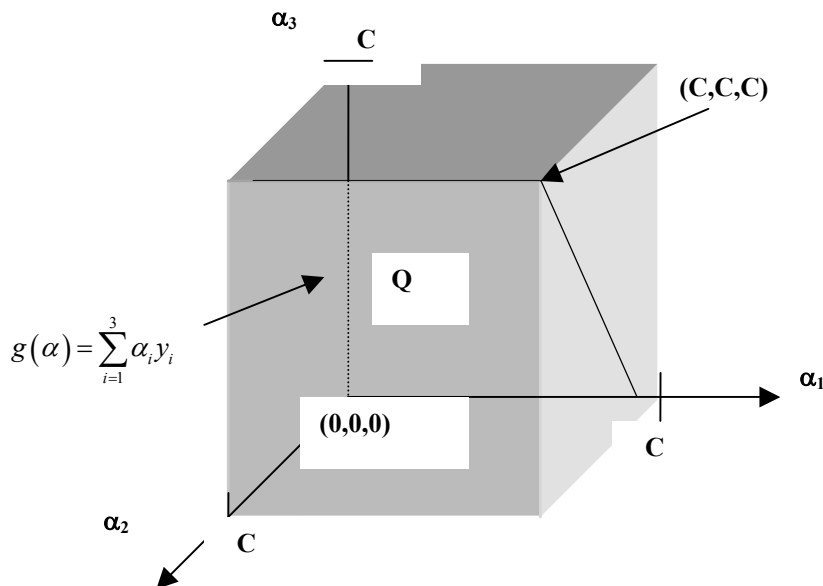


Fig IV.1: The feasible region  $\mathbb{Q}$ , for a classification problem with just 3 Lagrange Multipliers.

The problem with our now bounded region of solutions is that there could exist many linear manifolds (possibly infinite) all parallel to each other which gives a local minimum for P.3 and P.4. Optimization programs may then have slower rates of convergence if the solution set oscillates between these manifolds. In fact, Keerthi[16] has indirectly noted this problem with the Sequential Minimal Optimization(SMO) algorithm which he attributes to the pessimistic updates of the threshold,  $\mathbf{b}$  in the original SMO. If we look at the derivations of P.3 and P.4, we see that the linear restriction arises directly from treating the threshold  $\mathbf{b}$ , as a variable and is further evidence of our assertion that many possible linear manifolds exist which give equally good feasible regions. Adjusting the threshold  $\mathbf{b}$  during optimization, is then geometrically interpreted as alternating between the possible linear manifolds whilst looking for the optimal solution. Note that the position of the final separating hyperplane is also determined by the threshold  $\mathbf{b}$ .

There are a number of ways to avoid this problem. Firstly, we could remove the linear restriction completely from the quadratic program and end up with a feasible region which is a convex polyhedra due to H alone. We have discovered one way through the use of a special class of semiparametric functions which we have previously proposed in [14], but work on the generalization of this model has not yet been attempted. Others chose to incorporate the threshold as a Lagrange Multiplier instead[17] but this still requires explicit computation of  $\mathbf{b}$  in a different form. We have proposed several ways to eliminate the equality constraint. We may obtain a local solution instead but end up with an easier implementation with possibly better rates of convergence. However, further investigation is need to ascertain the effect of elimination on the generalization of the trained model.

### B. Dual Hyperplane Form

It is interesting to note that P.3 simplifies P.2 greatly by reducing the objective function to just one variable, that is the Lagrange Multiplier. In fact, what is happening is that we are now looking for the solution to  $\mathfrak{Z}(\mathbf{w}^*, \alpha)$  where the weight vector,  $\mathbf{w}^*$  is optimal provided  $\mathfrak{Z}(\mathbf{w}^*, \alpha^*)$  is a saddle point of  $\mathfrak{Z}$ . This is guaranteed anyway since  $\mathbf{w}^*$  is a saddle point obtained from the stationary gradient of P.2 in (10). Then by Kuhn and Tucker's saddle point theorem[18], we are assured that there exists an



optimal value for  $\alpha_0 \geq 0$  so that  $(\mathbf{w}^*, \alpha_0)$  is a saddle point for P.3 and hence also a solution to P.2. This allows us to assert that the dual hyperplane form (13) is derived only when  $\mathbf{w}$  is a saddle point for P.2 and hence implicitly embodies the minimization of the hyperplane margin. We note though that a modification of (13) has been used by Osuna[19] to approximate the decision surface, but his form can be shown to reduce to the original form by setting  $\beta_i = \alpha_i + \delta_i$  where  $\delta_i$  is some slack or tolerance variable. In our case, we argue that using this form in any objective function,  $\mathfrak{S}(w, \alpha)$  implies searching for the saddle point of  $\alpha$ , given that we have found the saddle point of  $\mathbf{w}$  or in short, maximized the margin of the hyperplane. Hence we conclude that the linear minimization programs solved by Joachims et. al. and also our proposed system of equations indirectly solve the quadratic Support Vector Machine program for classification.

## V. DISCUSSION

The derivation of programs P6-P8 is also indirectly related to solving a *linear complementarity problem*[8] which has been applied in some game theory and boundary value problems. Dantzig-Wolfe have used a principle pivoting method[9] which is an iterative method to find the solution to these kind of problems but in general all methods are closely related to solving a Linear Program and involves row operations on the equations in the objective function. Mangasarian[20, 21] further examines the solutions to linear complementarity programs if the objective function is concave. Concavity guarantees a solution which is difficult to prove for functions otherwise i.e nonconcave. However, the feasible region for our linear programs is bounded and it is known that there is a solution with every minimizer a relative boundary point and at least one minimizer an extreme point[22]. Geometrically, extreme points are vertices of the convex polyhedra which is the feasible region of solutions and if the unique minimizer to the unconstrained problem lies outside the polyhedra, the constrained minimizer lies on a boundary plane with at least one such minimizer being the vertex of the polyhedra.

The method of geometrical elimination which resulted in P8 interestingly gives rise to a number of possibilities for incremental SVM methods. Incremental methods which investigate different C parameters actually look at convex polyhedra of different sizes but the region of feasible solutions remain bounded to the surface of the hyperplane (42). The solution vector,  $\alpha$  of one C is now the known point  $\mathbf{p}$  for the next value of C. Based on this, it is possible to calculate notions of distances between optimal solutions using the distances between the points of the surface or even angles between different hyperplanes resulting from different normal vectors  $\mathbf{y}$  which arise due to additional input data.

The matrix equations of the linear program allows us to apply modified linear iterative methods and other methods that operate on sets of linear equations. This was meant to demonstrate a fast and simple way of implementing Support Vector Machines for classification. The single iterate updates gives us the possibility of using heuristics to choose the next point of update. We can cycle through the subsets of Support Vectors, arrange them in order of magnitudes or size of change and so on. Different datasets work well with different heuristics and we have yet to gain a full comprehension of the reasons for this observation. The convergence of this method has yet to be proved but we conjecture that it is similar to that of decomposition algorithms and activeset methods. The analysis of the rates of convergence of such methods will be the direction of our future work. In particular, we would be investigating the effect of ordering the elements in the working subsets on the convergence speeds.

## VI. CONCLUSION

We have transformed the quadratic program SVM classification problem to a set of linear equations which can be represented by a matrix equation solved on a compact set. Solving the system indirectly maximizes the margin of the hyperplane by using the dual hyperplane form. Future work would involve investigating this program applied to the Support Vector Regression model. Another interesting direction would be to investigate the effect of fixing the linear manifold on generalization accuracy. In fact, our previous work on semiparametric SVM classifiers can be interpreted as one example of fixing the linear manifold. The geometrical interpretation of the region of feasible solutions provides interesting ideas for the construction of optimization algorithms based on gradient flows in this bounded region.

REFERENCES

- [1] V. N. Vapnik, *The nature of statistical learning theory*, 2nd ed. New York: Springer, 2000.
- [2] B. Schölkopf and A. J. Smola, *Learning with kernels : support vector machines, regularization, optimization, and beyond*. Cambridge, Mass.: MIT Press, 2002.
- [3] C.-J. Lin, "LIBSVM," <http://www.csie.ntu.edu.tw/~cjlin/>.
- [4] 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.
- [5] 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.
- [6] 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.
- [7] B. Kristin and B. E., "Duality and Geometry in SVM Classifiers," presented at International Conference on Machine Learning, San Francisco, 2000.
- [8] R. Fletcher, *Practical methods of optimization*. Chichester [Eng.]; New York: J. Wiley, 1981.
- [9] P. E. Gill, W. Murray, and M. H. Wright, *Practical optimization*. London; New York: Academic Press, 1981.
- [10] J.-B. Hiriart-Urruty and C. Lemaréchal, *Convex analysis and minimization algorithms*. Berlin; New York: Springer-Verlag, 1993.
- [11] R. T. Rockafellar, *Convex analysis*. Princeton, N.J.: Princeton University Press, 1970.
- [12] C.-J. Lin, "A formal analysis of stopping criteria of decomposition methods for support vector machines," *Neural Networks, IEEE Transactions on*, vol. 13, pp. 1045-1052, 2002.
- [13] E. F. Osuna, R.; Girosit, F., "Training support vector machines: an application to face detection," presented at Computer Vision and Pattern Recognition, 1997. Proceedings., 1997 IEEE Computer Society Conference on, 1997.
- [14] D. Lai, M. Palaniswami, and N. Mani, "Fast Linear Optimization with Automatically Biased Support Vector Machines," Monash University MECSE-4-2003, 2003.
- [15] 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.
- [16] S. K. S. 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.
- [17] C.-J. Lin, "Linear convergence for a decomposition method for Support Vector Machines," November 2001.
- [18] J. Stoer and C. Witzgall, *Convexity and optimization in finite dimensions*. Berlin, New York,: Springer-Verlag, 1970.
- [19] B. Schölkopf, C. J. C. Burges, and A. J. Smola, *Advances in kernel methods : support vector learning*. Cambridge, Mass.: MIT Press, 1999.
- [20] O. L. Mangasarian, "Machine Learning via Polyhedral Concave Minimization," 95-20, 1995.
- [21] O. L. Mangasarian, "Solution of General Linear Complementarity Problems via Nondifferentiable Concave Minimization," 96-10, 1996.
- [22] W. Kaplan, *Maxima and minima with applications : practical optimization and duality*. New York: Wiley, 1999.