

Global Convergence of Decomposition Learning Methods for Support Vector Machines

Author(s): Norikazu Takahashi and Tetsuo Nishi

Journal: IEEE Transactions on Neural Networks

Volume: 17

Number: 6

Pages: 1362–1369

Month: November

Year: 2006

Published Version: http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=4012045

©2006 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

Global Convergence of Decomposition Learning Methods for Support Vector Machines

Norikazu Takahashi, *Member, IEEE* and Tetsuo Nishi, *Fellow, IEEE*

Abstract—Decomposition methods are well-known techniques for solving quadratic programming (QP) problems arising in support vector machines (SVMs). In each iteration of a decomposition method, a small number of variables are selected and a QP problem with only the selected variables is solved. Since large matrix computations are not required, decomposition methods are applicable to large QP problems. In this paper, we will make a rigorous analysis of the global convergence of general decomposition methods for SVMs. We first introduce a relaxed version of the optimality condition for the QP problems and then prove that a decomposition method reaches a solution satisfying this relaxed optimality condition within a finite number of iterations under a very mild condition on how to select variables.

Index Terms—Support vector machines, quadratic programming, decomposition method, global convergence, termination

I. INTRODUCTION

Support vector machines (SVMs) have recently attracted great attention in various fields such as pattern recognition, machine learning, neural networks, signal processing, and so on [1]–[4]. Given a set of l training samples, an SVM has to solve a quadratic programming (QP) problem with l variables. If the number of training samples are considerably large, the conventional QP solvers cannot be directly applied to SVM learning because large matrix computations are required. To overcome this difficulty, several techniques called decomposition methods have recently been proposed [5]–[8]. A basic strategy commonly used in decomposition methods is to execute two operations repeatedly until some optimality condition is satisfied; one is to select q variables among l and the other is to minimize the objective function by updating only the the selected q variables. The set of q variables selected for updating is called *the working set*. An example of decomposition methods is the sequential minimal optimization (SMO) algorithm proposed by Platt [5] in which only two variables are selected for the working set in each iteration. Another example is SVM^{light} [6], the most widely used learning algorithm for SVMs, in which q , the size of the working set, can be set to any even number. Since large matrix computations are not required, decomposition methods are

useful for SVMs especially when a large number of training samples are given.

For any decomposition method, it is very important to guarantee that the sequence of solutions converges to an optimal solution within a finite number of iterations. Some theoretical results concerning the global convergence of decomposition methods can be found in the literature [9]–[16]. Keerthi and Gilbert [9] analyzed the behavior of the generalized SMO algorithm and gave a proof that the algorithm terminates within a finite number of iterations under a pre-specified stopping condition and tolerance¹. Lin [10] proved that the sequence of solutions obtained by an SMO algorithm converges asymptotically to an optimal solution. Recently, Chen *et al.* [12] carried out a comprehensive study on SMO algorithms and gave several results on asymptotic convergence, finite termination, shrinking and caching, and convergence rate in a general setting on the working set selection. However, the results given in [9]–[12] rely heavily on the fact that only two variables are updated in each iteration, and therefore it seems difficult to extend to the general case where q is greater than two.

Chang *et al.* [13] considered a kind of decomposition method in which q is not restricted to two, and proved the convergence of the method. However, the working set selection in their method is formulated in a rather special form, and therefore the result does not hold for other decomposition methods. Lin [14] investigated the properties of the sequence of solutions obtained by SVM^{light} in detail and proved that any limit point of the sequence is an optimal solution. Lin [15] also showed that a class of decomposition methods including SVM^{light} stops within a finite number of iterations if a relaxed version of stopping criterion is used. These results are important because the convergence of SVM^{light} was guaranteed theoretically for the first time. However, on the other hand, these results require a stronger assumption on the Hessian matrix of the objective function than positive semi-definiteness [14, Assumption IV.1], except the special case where $q = 2$ [10]. It is known that this assumption is satisfied if the Gaussian kernel is used and training samples are different from each other, but, it is also known that this assumption does not hold true for any kernel function if training samples contain two or more identical data [10]. List and Simon [16] considered a general class of QP problems which includes the one arising in SVMs, and shown that a decomposition method converges to an optimal solution if its working set selection method satisfies three abstract conditions. This result is general in the sense that no specific working set selection method is

This work was supported in part by the 21st Century COE (Center of Excellence) Program “Reconstruction of Social Infrastructure Related to Information Science and Electrical Engineering”.

N. Takahashi is with the Department of Computer Science and Communication Engineering, Kyushu University, Fukuoka, 812-8581 Japan (e-mail: norikazu@csce.kyushu-u.ac.jp).

T. Nishi is with the Department of Computer Science and Communication Engineering, Kyushu University, Fukuoka, 812-8581 Japan. He is now with the Faculty of Science and Engineering, Waseda University, Tokyo, 169-0072 Japan (e-mail: nishi-t@waseda.jp).

¹The authors of the present paper have recently pointed out the incompleteness of their proof and given a more rigorous proof [11].

assumed, but still requires the same assumption on the Hessian matrix of the objective function as [14] and [15].

In this paper, we will present a new convergence proof for general decomposition methods for SVMs. Unlike the existing results mentioned above, we will neither restrict ourselves to a specific working set selection nor assume any condition on the Hessian matrix of the objective function except for positive semi-definiteness. Instead, we will employ a relaxed version of the Karush-Kuhn-Tucker (KKT) conditions as the optimality condition. The relaxed version contains two positive parameters τ and ϵ , and approaches the strict KKT conditions as both τ and ϵ go to zero. By using the global convergence theorem in optimization theory [17], we will prove that the decomposition method stops within a finite number of iterations after finding an optimal solution for any τ and ϵ if the working set contains at least one pair of indices violating the optimality condition in each iteration. This condition on the working set selection is very mild and thus can be applied to many decomposition methods.

In our convergence proof, closedness of a point-to-set map [17] plays a central role. It is important to note that the point-to-set map defined for the decomposition method considered in this paper is not closed if either the strict KKT condition or the relaxed version which has been often used so far [6], [9], [11], [15] is employed as the optimality condition. Indeed this is the main difficulty on proving the convergence in earlier works. An extensive discussion is given in [14].

This paper is organized as follows. In Section II, the global convergence theorem for general optimization problems and some related results are reviewed for later discussions. In Section III, the QP problem arising in SVMs, two types of optimality conditions, and the algorithm of general decomposition methods are explained. In Section IV, the relationship between two optimality conditions is firstly discussed and then convergence theorems for decomposition methods are proved. In Section V, these convergence theorems are applied to some well-known learning algorithms for SVMs. Finally, concluding remarks are given in Section VI.

II. PRELIMINARIES

First of all, we will review a fundamental result in optimization theory known as the global convergence theorem [17]. This theorem plays a central role in our global convergence analysis of decomposition methods.

Let us consider the following optimization problem.

Problem 1: Find \mathbf{x} which minimizes the objective function $f(\mathbf{x})$ under the constraint $\mathbf{x} \in X$.

An algorithm for solving Problem 1 can be viewed as an iterative process that generates a sequence $\{\mathbf{x}(k)\}_{k=0}^{\infty}$ by

$$\mathbf{x}(k+1) \in \mathbf{A}(\mathbf{x}(k)), \quad k = 0, 1, \dots \quad (1)$$

where $\mathbf{x}(0) \in X$ is a given initial point and \mathbf{A} is a point-to-set map that assigns to each point in the domain X a subset of X . According to this definition, we can state that developing an algorithm for solving Problem 1 is equivalent to determining the point-to-set map \mathbf{A} . Apparently, the most desirable property of the map \mathbf{A} is that the sequence $\{\mathbf{x}(k)\}_{k=0}^{\infty}$ generated

by (1) converges to an optimal solution of Problem 1 for any initial point $\mathbf{x}(0) \in X$.

Definition 1: Let X , Ω and \mathbf{A} be a nonempty closed set in \mathbb{R}^n , a subset of X , and a point-to-set map from X to X , respectively. If a continuous function Z defined in X satisfies

$$\text{if } \mathbf{x} \notin \Omega \text{ and } \mathbf{y} \in \mathbf{A}(\mathbf{x}) \text{ then } Z(\mathbf{y}) < Z(\mathbf{x})$$

then Z is called a *descent function* for Ω and \mathbf{A} .

Definition 2: Let X and Y be nonempty closed sets in \mathbb{R}^n and \mathbb{R}^m , respectively. Let \mathbf{A} be a point-to-set map from X to Y . The map \mathbf{A} is said to be *closed* at $\bar{\mathbf{x}}$ if the two assumptions

- 1) $\mathbf{x}(k) \in X$, $\forall k$ and $\lim_{k \rightarrow \infty} \mathbf{x}(k) = \bar{\mathbf{x}}$
- 2) $\mathbf{y}(k) \in \mathbf{A}(\mathbf{x}(k))$, $\forall k$ and $\lim_{k \rightarrow \infty} \mathbf{y}(k) = \bar{\mathbf{y}}$

imply that $\bar{\mathbf{y}} \in \mathbf{A}(\bar{\mathbf{x}})$. The map \mathbf{A} is said to be closed on X if it is closed at each point in X .

In terms of these definitions, the global convergence theorem can be stated as follows.

Theorem 1 ([17]): Let X , Ω and \mathbf{A} be a nonempty closed set in \mathbb{R}^n , a subset of X , and a point-to-set map from X to X , respectively. Let $\{\mathbf{x}(k)\}_{k=0}^{\infty}$ be a sequence generated by (1) with $\mathbf{x}(0) \in X$. Every convergent subsequence of $\{\mathbf{x}(k)\}_{k=0}^{\infty}$ has a limit in Ω if the following conditions are satisfied.

- 1) For all k , $\mathbf{x}(k)$ belongs to a compact set $S \subseteq X$.
- 2) There exists a descent function Z for Ω and \mathbf{A} .
- 3) The map \mathbf{A} is closed on $X \setminus \Omega$.

Note that Ω in Theorem 1 can be any subset of X . In particular, if Ω is set to the set of optimal solutions of Problem 1, Theorem 1 gives a sufficient condition for the sequence $\{\mathbf{x}(k)\}_{k=0}^{\infty}$ generated by (1) to converge to an optimal solution. Theorem 1 is hence a very useful tool for proving the global convergence of an algorithm. However, it is difficult in general to develop a map \mathbf{A} having the global convergence property for the set of optimal solutions. Thus Ω is often set to, for example, $\{\mathbf{x} \mid \mathbf{x} \text{ is a local optimal solution}\}$ and $\{\mathbf{x} \mid \mathbf{x} \in X, f(\mathbf{x}) \leq b\}$ where b is a constant.

The following lemma can easily be obtained.

Lemma 1: Let $\mathbf{A} : X \rightarrow Y$ and $\mathbf{B} : X \rightarrow Y$ be point-to-set maps. If both \mathbf{A} and \mathbf{B} are closed at \mathbf{x} , then the map $\mathbf{C}(\mathbf{x}) = \mathbf{A}(\mathbf{x}) \cup \mathbf{B}(\mathbf{x})$ is also closed at \mathbf{x} .

Proof: Let $\{\mathbf{x}(k)\}_{k=0}^{\infty}$ be any sequence such that $\mathbf{x}(k) \in X$, $\forall k$ and $\lim_{k \rightarrow \infty} \mathbf{x}(k) = \mathbf{x}$. Let $\{\mathbf{y}(k)\}_{k=0}^{\infty}$ be any sequence such that $\mathbf{y}(k) \in \mathbf{C}(\mathbf{x}(k))$, $\forall k$ and $\lim_{k \rightarrow \infty} \mathbf{y}(k) = \mathbf{y}$. Then it is obvious that at least one of two statements:

- 1) $\mathbf{y}(k) \in \mathbf{A}(\mathbf{x}(k))$ for infinitely many k ,
- 2) $\mathbf{y}(k) \in \mathbf{B}(\mathbf{x}(k))$ for infinitely many k

holds. We will assume without loss of generality that the first case holds true. Let the set of all k such that $\mathbf{y}(k) \in \mathbf{A}(\mathbf{x}(k))$ be denoted by K_A . Then the sequence $\{\mathbf{x}(k)\}_{k \in K_A}$ satisfies $\mathbf{x}(k) \in X$, $\forall k \in K_A$ and $\lim_{k \rightarrow \infty, k \in K_A} \mathbf{x}(k) = \mathbf{x}$. Also, the sequence $\{\mathbf{y}(k)\}_{k \in K_A}$ satisfies $\mathbf{y}(k) \in \mathbf{A}(\mathbf{x}(k))$, $\forall k \in K_A$ and $\lim_{k \rightarrow \infty, k \in K_A} \mathbf{y}(k) = \mathbf{y}$. Since the point-to-set map \mathbf{A} is closed at \mathbf{x} , we have $\mathbf{y} \in \mathbf{A}(\mathbf{x})$ which implies that $\mathbf{y} \in \mathbf{C}(\mathbf{x})$. Therefore the point-to-set map \mathbf{C} is closed at \mathbf{x} . ■

A result concerning the closedness of composite maps will also be needed in later discussions.

Definition 3: Let $\mathbf{A} : X \rightarrow Y$ and $\mathbf{B} : Y \rightarrow Z$ be point-to-set maps. The composite map $\mathbf{C} = \mathbf{B}\mathbf{A}$ is defined as the

point-to-set map $C : X \rightarrow Z$ with

$$C(\mathbf{x}) = \cup_{\mathbf{y} \in A(\mathbf{x})} B(\mathbf{y}).$$

Lemma 2 ([17]): Let $A : X \rightarrow Y$ and $B : Y \rightarrow Z$ be point-to-set maps. If A is closed at \mathbf{x} , B is closed on $A(\mathbf{x})$ and Y is compact, then the composite map $C = BA$ is closed at \mathbf{x} .

III. DECOMPOSITION METHOD

A. SVM Dual Problem

Suppose that we are given a set of l training samples $\{(\mathbf{p}_i, d_i)\}_{i=1}^l$ where $\mathbf{p}_i \in \mathbb{R}^n$ is the i -th input pattern and $d_i \in \{1, -1\}$ represents the class to which \mathbf{p}_i belongs. The learning of an SVM with the kernel function $K(\cdot, \cdot)$ leads to the following QP problem (for more details on formulation, see for example [1]).

Problem 2: Find $\boldsymbol{\alpha} = [\alpha_1, \alpha_2, \dots, \alpha_l]^T$ which minimizes the objective function

$$W(\boldsymbol{\alpha}) = \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l q_{ij} \alpha_i \alpha_j - \sum_{i=1}^l \alpha_i$$

under the constraints

$$\sum_{i=1}^l d_i \alpha_i = 0 \quad (2)$$

$$0 \leq \alpha_i \leq C, \quad i = 1, 2, \dots, l \quad (3)$$

where $q_{ij} = d_i d_j K(\mathbf{p}_i, \mathbf{p}_j)$ and C is a user-specified positive constant.

Throughout this paper, we assume that the kernel function $K(\cdot, \cdot)$ satisfies Mercer's condition [1]. In this case, $\mathbf{Q} = [q_{ij}] \in \mathbb{R}^{l \times l}$ is a positive semi-definite matrix, and therefore Problem 2 is a convex QP problem. Note that the optimal solution of Problem 2 is not necessarily unique because $W(\boldsymbol{\alpha})$ is not strictly convex. The feasible region of Problem 2, that is, the set of $\boldsymbol{\alpha}$ satisfying (2) and (3), is denoted by S which is apparently a compact set. Also, the set $\{1, 2, \dots, l\}$ is denoted by L .

B. Optimality Conditions

Since Problem 2 is a convex QP problem, optimal solutions are completely characterized by the KKT conditions [17], that is, $\boldsymbol{\alpha} \in S$ is an optimal solution of Problem 2 if and only if there exist constants $\lambda, \mu_1, \mu_2, \dots, \mu_l, \nu_1, \nu_2, \dots, \nu_l$ such that

$$\begin{cases} \partial W(\boldsymbol{\alpha}) / \partial \alpha_i + \lambda d_i - \mu_i + \nu_i = 0 \\ \mu_i \alpha_i = 0 \\ \nu_i (\alpha_i - C) = 0 \\ \mu_i \geq 0 \\ \nu_i \geq 0 \end{cases} \quad (4)$$

for all $i \in L$. As shown in [9] and [15], this condition can be rewritten in a more compact form as

$$\min_{i \in I_{\text{up}}(\boldsymbol{\alpha})} F_i(\boldsymbol{\alpha}) \geq \max_{i \in I_{\text{low}}(\boldsymbol{\alpha})} F_i(\boldsymbol{\alpha}) \quad (5)$$

where

$$F_i(\boldsymbol{\alpha}) = d_i \left(\sum_{j=1}^l q_{ij} \alpha_j - 1 \right),$$

$$I_{\text{up}}(\boldsymbol{\alpha}) = \{i \mid \alpha_i < C, d_i = 1\} \cup \{i \mid \alpha_i > 0, d_i = -1\},$$

$$I_{\text{low}}(\boldsymbol{\alpha}) = \{i \mid \alpha_i < C, d_i = -1\} \cup \{i \mid \alpha_i > 0, d_i = 1\}.$$

Let Ω^* denote the set of optimal solutions of Problem 2. Then, by using the above notations, we can express Ω^* as follows:

$$\Omega^* = \{\boldsymbol{\alpha} \in S \mid \min_{i \in I_{\text{up}}(\boldsymbol{\alpha})} F_i(\boldsymbol{\alpha}) \geq \max_{i \in I_{\text{low}}(\boldsymbol{\alpha})} F_i(\boldsymbol{\alpha})\}.$$

In a practical situation, the optimality condition (5) is often relaxed as

$$\min_{i \in I_{\text{up}}(\boldsymbol{\alpha})} F_i(\boldsymbol{\alpha}) \geq \max_{i \in I_{\text{low}}(\boldsymbol{\alpha})} F_i(\boldsymbol{\alpha}) - \tau \quad (6)$$

where τ is a positive constant [6], [7], [15].

In this paper, on the other hand, we employ neither (5) nor (6) but the inequality:

$$\min_{i \in I_{\text{up}}^\epsilon(\boldsymbol{\alpha})} F_i(\boldsymbol{\alpha}) > \max_{i \in I_{\text{low}}^\epsilon(\boldsymbol{\alpha})} F_i(\boldsymbol{\alpha}) - \tau \quad (7)$$

for the optimality condition, where

$$I_{\text{up}}^\epsilon(\boldsymbol{\alpha}) = \{i \mid \alpha_i \leq C - \epsilon, d_i = 1\} \cup \{i \mid \alpha_i \geq \epsilon, d_i = -1\},$$

$$I_{\text{low}}^\epsilon(\boldsymbol{\alpha}) = \{i \mid \alpha_i \leq C - \epsilon, d_i = -1\} \cup \{i \mid \alpha_i \geq \epsilon, d_i = 1\},$$

and ϵ is any positive constant smaller than $C/2$. Usually ϵ is set to a sufficiently small positive number. The role of $I_{\text{up}}^\epsilon(\boldsymbol{\alpha})$ and $I_{\text{low}}^\epsilon(\boldsymbol{\alpha})$ is to consider α_i which is sufficiently close to 0 (C , resp.) to be exactly 0 (C , resp.). This is a technique used in the implementation of SVM^{light} algorithm [6] (see the source code² of SVM^{light} developed by Joachims). In the following, any $\boldsymbol{\alpha} \in S$ satisfying (7) is said to be a (τ, ϵ) -optimal solution. The set of (τ, ϵ) -optimal solutions is denoted by $\Omega^{(\tau, \epsilon)}$, that is,

$$\Omega^{(\tau, \epsilon)} = \{\boldsymbol{\alpha} \in S \mid \min_{i \in I_{\text{up}}^\epsilon(\boldsymbol{\alpha})} F_i(\boldsymbol{\alpha}) > \max_{i \in I_{\text{low}}^\epsilon(\boldsymbol{\alpha})} F_i(\boldsymbol{\alpha}) - \tau\}.$$

Also, a pair of indices (i, j) such that

$$i \in I_{\text{up}}^\epsilon(\boldsymbol{\alpha}), j \in I_{\text{low}}^\epsilon(\boldsymbol{\alpha}), F_i(\boldsymbol{\alpha}) \leq F_j(\boldsymbol{\alpha}) - \tau$$

is called a (τ, ϵ) -violating pair at $\boldsymbol{\alpha}$. If a pair (i, j) is not a (τ, ϵ) -violating pair at $\boldsymbol{\alpha}$, the pair is called a (τ, ϵ) -feasible pair at $\boldsymbol{\alpha}$. It is obvious from these definitions that $\boldsymbol{\alpha} \in S$ is a (τ, ϵ) -optimal solution if and only if there is no (τ, ϵ) -violating pair at $\boldsymbol{\alpha}$.

Properties of the two sets Ω^* and $\Omega^{(\tau, \epsilon)}$ as well as the relationship between them will be studied in detail in the next section.

C. Algorithm of Decomposition Method

A basic strategy commonly used in all decomposition methods is to repeat two operations until some optimality condition is satisfied; one is to select q variables among l for the working set and the other is to minimize the objective function $W(\boldsymbol{\alpha})$ by updating only the selected q variables. This is formally expressed as follows:

²Source code and binaries of SVM^{light} are available at <http://svmlight.joachims.org/>.

Algorithm 1: Given training samples $\{(\mathbf{p}_i, d_i)\}_{i=1}^l$, a kernel function $K(\cdot, \cdot)$, a positive constant C and an integer $q(\leq l)$, execute the following procedures.

- 1) Let $\alpha(0) = \mathbf{0}$ and $k = 0$.
- 2) If $\alpha = \alpha(k)$ satisfies the optimality condition (7) then stop. Otherwise go to Step 3).
- 3) Select the working set $L_B(k) \subseteq L = \{1, 2, \dots, l\}$ where $|L_B(k)| \leq q$.
- 4) Find $\alpha = [\alpha_1, \alpha_2, \dots, \alpha_l]^T$ which minimizes the objective function $W(\alpha)$ under the constraints (2), (3) and $\alpha_i = \alpha_i(k), \forall i \in L_N(k) = L \setminus L_B(k)$.
- 5) Set $\alpha(k+1)$ to an optimal solution of the optimization problem in Step 4). Add 1 to k , and go to Step 2).

We will refer to the optimization problem in Step 4) as the subproblem in the following. Since the subproblem has at most q variables, the amount of memory required for Algorithm 1 is linear in l , while l^2 elements of the matrix \mathbf{Q} must be stored in memory if Problem 2 is solved at a time. This is the main advantage of decomposition methods. As for computation time, each subproblem can be solved faster as q decreases. In particular, in the case where $q = 2$, subproblems can be solved analytically and hence Algorithm 1 can be implemented without QP solvers [5]. However, it should be noted that the number of iterations increases as q decreases in general.

It is apparent that the sequence $\{\alpha(k)\}_{k=1}^\infty$ generated by Algorithm 1 satisfies two conditions:

$$\alpha(k) \in S \quad (8)$$

$$W(\alpha(k+1)) \leq W(\alpha(k)) \quad (9)$$

for all k . Since the objective function $W(\cdot)$ is bounded from below in S , Eq.(9) implies that the sequence $\{W(\alpha(k))\}_{k=0}^\infty$ necessarily converges to a certain value. However, on the other hand, it is not clear whether the sequence $\{\alpha(k)\}_{k=0}^\infty$ converges to $\Omega^{(\tau, \epsilon)}$ or not.

The optimality condition for the subproblem is expressed in the same form as (5), that is, for given $L_B(k) \subseteq L$ and $\alpha(k) \in S$, $\alpha \in S$ is an optimal solution of the subproblem if and only if the following conditions are satisfied.

$$\begin{cases} \min_{i \in I_{\text{up}}(\alpha) \cap L_B(k)} F_i(\alpha) \geq \max_{i \in I_{\text{low}}(\alpha) \cap L_B(k)} F_i(\alpha) \\ \alpha_i = \alpha_i(k), \quad \forall i \in L_N(k) = L \setminus L_B(k) \end{cases}$$

IV. GLOBAL CONVERGENCE ANALYSIS OF DECOMPOSITION METHODS

A. Properties of Ω^* and $\Omega^{(\tau, \epsilon)}$

Before proceeding to the convergence analysis of Algorithm 1, we study properties of Ω^* and $\Omega^{(\tau, \epsilon)}$.

Lemma 3: $I_{\text{up}}(\alpha) \supseteq I_{\text{up}}^\epsilon(\alpha)$ and $I_{\text{low}}(\alpha) \supseteq I_{\text{low}}^\epsilon(\alpha)$ for any $\alpha \in S$ and $\epsilon \in (0, C/2)$.

Proof: We will prove only the first formula. The second one can be proved similarly. Let i be any member of $I_{\text{up}}^\epsilon(\alpha)$. Then α_i satisfies $0 \leq \alpha_i < C - \epsilon$ if $d_i = 1$ and $\epsilon < \alpha_i \leq C$ if $d_i = -1$. In the former case, i belongs to $I_{\text{up}}(\alpha)$ because either $0 < \alpha_i < C$ or $\alpha_i = 0$ holds. In the latter case, i belongs to $I_{\text{up}}(\alpha)$ because either $0 < \alpha_i < C$ or $\alpha_i = C$ holds. Therefore, any member of $I_{\text{up}}^\epsilon(\alpha)$ belongs to $I_{\text{up}}(\alpha)$. This implies $I_{\text{up}}(\alpha) \supseteq I_{\text{up}}^\epsilon(\alpha)$. ■

Proposition 1: $\Omega^{(\tau, \epsilon)} \supseteq \Omega^*$ for any $\tau > 0$ and $\epsilon \in (0, C/2)$.

Proof: Let α be any point in Ω^* . Then α satisfies (5). It follows from Lemma 3 that

$$\min_{i \in I_{\text{up}}^\epsilon(\alpha)} F_i(\alpha) \geq \min_{i \in I_{\text{up}}(\alpha)} F_i(\alpha),$$

$$\max_{i \in I_{\text{low}}(\alpha)} F_i(\alpha) \geq \max_{i \in I_{\text{low}}^\epsilon(\alpha)} F_i(\alpha).$$

From these two inequalities and (5), we have

$$\min_{i \in I_{\text{up}}^\epsilon(\alpha)} F_i(\alpha) \geq \max_{i \in I_{\text{low}}(\alpha)} F_i(\alpha)$$

which implies $\alpha \in \Omega^{(\tau, \epsilon)}$. ■

Lemma 4: Let $\{\alpha(n)\}_{n=0}^\infty$ be any sequence such that $\alpha(n) \in S, \forall n$ and $\lim_{n \rightarrow \infty} \alpha(n) = \bar{\alpha}$. Then there exist positive integers n_1 and n_2 such that

$$I_{\text{up}}(\alpha(n)) \supseteq I_{\text{up}}(\bar{\alpha}), \quad \forall n \geq n_1$$

$$I_{\text{low}}(\alpha(n)) \supseteq I_{\text{low}}(\bar{\alpha}), \quad \forall n \geq n_2$$

Proof: We will prove only the first formula. The second one can be proved in the same way. Let i be any member of $I_{\text{up}}(\bar{\alpha})$. Then $\bar{\alpha}_i$ satisfies $\bar{\alpha}_i < C$ if $d_i = 1$ and $\bar{\alpha}_i > 0$ if $d_i = -1$. In the former case, since $\alpha_i(n)$ converges to $\bar{\alpha}_i$, there exists a positive integer $n_1(i)$ such that $\alpha_i(n) < C, \forall n \geq n_1(i)$ which implies $i \in I_{\text{up}}(\alpha(n)), \forall n \geq n_1(i)$. In the latter case, it is shown in the same way that there exists a positive integer $n_1(i)$ such that $i \in I_{\text{up}}(\alpha(n)), \forall n \geq n_1(i)$. Let $n_1 = \max_{i \in I_{\text{up}}(\bar{\alpha})} n_1(i)$. Then all members of $I_{\text{up}}(\bar{\alpha})$ belong to $I_{\text{up}}(\alpha(n)), \forall n \geq n_1$. This completes the proof. ■

This lemma was first given by Lin [14, Lemma IV.4]. We have just restated the result in terms of our notations, and given a proof for the sake of the reader's convenience.

Proposition 2: The set Ω^* is closed.

Proof: Let $\{\alpha(n)\}_{n=1}^\infty$ be any sequence such that $\alpha(n) \in \Omega^*, \forall n$ and $\lim_{n \rightarrow \infty} \alpha(n) = \bar{\alpha}$. It suffices for us to show that $\bar{\alpha} \in \Omega^*$. Since $\alpha(n) \in \Omega^*, \forall n$, we have

$$\min_{i \in I_{\text{up}}(\alpha(n))} F_i(\alpha(n)) \geq \max_{i \in I_{\text{low}}(\alpha(n))} F_i(\alpha(n)), \quad \forall n.$$

It follows from this inequality and Lemma 4 that there exists a positive integer n_1 such that

$$\min_{i \in I_{\text{up}}(\bar{\alpha})} F_i(\alpha(n)) \geq \max_{i \in I_{\text{low}}(\bar{\alpha})} F_i(\alpha(n)), \quad \forall n \geq n_1.$$

By letting n go to infinity, we have

$$\min_{i \in I_{\text{up}}(\bar{\alpha})} F_i(\bar{\alpha}) \geq \max_{i \in I_{\text{low}}(\bar{\alpha})} F_i(\bar{\alpha})$$

which implies $\bar{\alpha} \in \Omega^*$. Therefore, Ω^* is a closed set. ■

Lemma 5: Let $\{\alpha(n)\}_{n=0}^\infty$ be any sequence such that $\alpha(n) \in S, \forall n$ and $\lim_{n \rightarrow \infty} \alpha(n) = \bar{\alpha}$. Then there exist positive integers n_1 and n_2 such that

$$I_{\text{up}}^\epsilon(\alpha(n)) \subseteq I_{\text{up}}^\epsilon(\bar{\alpha}), \quad \forall n \geq n_1$$

$$I_{\text{low}}^\epsilon(\alpha(n)) \subseteq I_{\text{low}}^\epsilon(\bar{\alpha}), \quad \forall n \geq n_2$$

for any $\epsilon \in (0, C/2)$.

Proof: We will prove only the first formula. The second one can be proved similarly. Let i be any nonmember of $I_{\text{up}}^\epsilon(\bar{\alpha})$. Then $\bar{\alpha}_i$ satisfies $\bar{\alpha}_i > C - \epsilon$ if $d_i = 1$ and $\bar{\alpha}_i < \epsilon$

if $d_i = -1$. In the former case, since $\alpha_i(n)$ converges to $\bar{\alpha}_i$, there exists a positive integer $n_1(i)$ such that $\alpha_i(n) > C - \epsilon$, $\forall n \geq n_1(i)$ which implies $i \notin I_{\text{up}}^\epsilon(\alpha(n))$, $\forall n \geq n_1(i)$. In the latter case, it is shown in the same way that there exists a positive integer $n_1(i)$ such that $i \notin I_{\text{up}}^\epsilon(\alpha(n))$, $\forall n \geq n_1(i)$. Let $n_1 = \max_{i \notin I_{\text{up}}^\epsilon(\bar{\alpha})} n_1(i)$. Then all nonmembers of $I_{\text{up}}(\bar{\alpha})$ do not belong to $I_{\text{up}}^\epsilon(\alpha(n))$, $\forall n \geq n_1$. This is equivalent to the first formula. ■

Proposition 3: The set $S \setminus \Omega^{(\tau, \epsilon)}$ is closed for any $\tau > 0$ and $\epsilon \in (0, C/2)$.

Proof: Let $\{\alpha(n)\}_{n=1}^\infty$ be any sequence such that $\alpha(n) \in S \setminus \Omega^{(\tau, \epsilon)}$, $\forall n$ and $\lim_{n \rightarrow \infty} \alpha(n) = \bar{\alpha}$. Then we have

$$\min_{i \in I_{\text{up}}^\epsilon(\alpha(n))} F_i(\alpha(n)) \leq \max_{i \in I_{\text{low}}^\epsilon(\alpha(n))} F_i(\alpha(n)) - \tau, \forall n.$$

It follows from this inequality and Lemma 5 that there exists a positive integer n_1 such that

$$\min_{i \in I_{\text{up}}^\epsilon(\bar{\alpha})} F_i(\bar{\alpha}) \leq \max_{i \in I_{\text{low}}^\epsilon(\bar{\alpha})} F_i(\bar{\alpha}) - \tau, \forall n \geq n_1.$$

By letting n go to infinity in both sides, we have

$$\min_{i \in I_{\text{up}}^\epsilon(\bar{\alpha})} F_i(\bar{\alpha}) \leq \max_{i \in I_{\text{low}}^\epsilon(\bar{\alpha})} F_i(\bar{\alpha}) - \tau$$

which means $\bar{\alpha} \in S \setminus \Omega^{(\tau, \epsilon)}$. Therefore, $\Omega^{(\tau, \epsilon)}$ is a closed set for any $\tau > 0$ and $\epsilon \in (0, C/2)$. ■

Proposition 4: The set $\Omega^{(\tau, \epsilon)}$ converges to Ω^* as the positive constants τ and ϵ approach 0.

Proof: One can easily see that $\lim_{\epsilon \rightarrow 0+} I_{\text{up}}^\epsilon(\alpha) = I_{\text{up}}(\alpha)$ and $\lim_{\epsilon \rightarrow 0+} I_{\text{low}}^\epsilon(\alpha) = I_{\text{low}}(\alpha)$, where $\epsilon \rightarrow 0+$ means ϵ approaches 0 from right. Thus we have

$$\begin{aligned} \lim_{\epsilon \rightarrow 0+} \Omega^{(\tau, \epsilon)} \\ = \{ \alpha \in S \mid \min_{i \in I_{\text{up}}(\alpha)} F_i(\alpha) > \max_{i \in I_{\text{low}}(\alpha)} F_i(\alpha) - \tau \}. \end{aligned}$$

Furthermore, if τ approaches 0 from right, then the right-hand side of the above inequality converges to Ω^* , that is, $\lim_{\tau \rightarrow 0+} \lim_{\epsilon \rightarrow 0+} \Omega^{(\tau, \epsilon)} = \Omega^*$. ■

B. Convergence Proof

Let $V_q(\alpha)$ be the family of sets $M \subseteq L$ such that $|M| \leq q$ and M contains at least one (τ, ϵ) -violating pair at $\alpha \in S$. Then the following lemma holds.

Lemma 6: Let $\{\alpha(n)\}_{n=0}^\infty$ be any sequence such that $\alpha(n) \in S$, $\forall n$ and $\lim_{n \rightarrow \infty} \alpha(n) = \bar{\alpha}$. If $\bar{\alpha} \in S \setminus \Omega^{(\tau, \epsilon)}$ then

$$V_q(\alpha(n)) \subseteq V_q(\bar{\alpha}) \quad (10)$$

for sufficiently large n .

Proof: Let (i, j) be any (τ, ϵ) -feasible pair at $\bar{\alpha}$. Then at least one of the following three conditions holds.

- 1) $i \notin I_{\text{up}}^\epsilon(\bar{\alpha})$
- 2) $j \notin I_{\text{low}}^\epsilon(\bar{\alpha})$
- 3) $i \in I_{\text{up}}^\epsilon(\bar{\alpha})$, $j \in I_{\text{low}}^\epsilon(\bar{\alpha})$, $F_i(\bar{\alpha}) > F_j(\bar{\alpha}) - \tau$

In Case 1), it is easily seen from Lemma 5 that $i \notin I_{\text{up}}^\epsilon(\alpha(n))$ for sufficiently large n , which means (i, j) is a (τ, ϵ) -feasible pair at $\alpha(n)$ for sufficiently large n . In Case 2), we can draw

the same conclusion as Case 1). In Case 3), it follows from Lemma 5 and the continuity of $F_i(\cdot)$ that

$$i \in I_{\text{up}}^\epsilon(\alpha(n)), j \in I_{\text{low}}^\epsilon(\alpha(n)), F_i(\alpha(n)) > F_j(\alpha(n)) - \tau$$

holds for sufficiently large n . This means (i, j) is a (τ, ϵ) -feasible pair at $\alpha(n)$ for sufficiently large n . Therefore, in all cases, the set of (τ, ϵ) -feasible pairs at $\bar{\alpha}$ is included in that at $\alpha(n)$ for sufficiently large n . Conversely, the set of (τ, ϵ) -violating pairs at $\alpha(n)$ is included in that at $\bar{\alpha}$ for sufficiently large n . Eq.(10) is immediately derived from this fact and the definition of $V_q(\cdot)$. ■

For any $M \subseteq L$ and $\alpha \in S$, we define the point-to-set map $\Gamma_M(\alpha)$ as

$$\begin{aligned} \Gamma_M(\alpha) \triangleq \{ \mathbf{y} \in S \mid y_i = \alpha_i \ \forall i \in L \setminus M, \\ \min_{i \in I_{\text{up}}(\mathbf{y}) \cap M} F_i(\mathbf{y}) \geq \max_{i \in I_{\text{low}}(\mathbf{y}) \cap M} F_i(\mathbf{y}) \}. \end{aligned}$$

By using this definition, the set of optimal solutions of the subproblem in Step 4) can be expressed as $\Gamma_{L_B(k)}(\alpha(k))$. We also define a point-to-set map \mathbf{A} from S to itself as follows:

$$\mathbf{A}(\alpha) = \begin{cases} \cup_{M \in V_q(\alpha)} \Gamma_M(\alpha), & \text{if } \alpha \notin \Omega^{(\tau, \epsilon)} \\ \alpha, & \text{if } \alpha \in \Omega^{(\tau, \epsilon)} \end{cases} \quad (11)$$

Lemma 7: For any $M \subseteq L$, the point-to-set map $\Gamma_M(\alpha)$ is closed on S .

Proof: Let $\{\alpha(n)\}_{n=0}^\infty$ be any sequence such that $\alpha(n) \in S$, $\forall n$ and $\lim_{n \rightarrow \infty} \alpha(n) = \bar{\alpha} \in S$. Let $\{\beta(n)\}_{n=0}^\infty$ be any sequence such that $\beta(n) \in \Gamma_M(\alpha(n))$, $\forall n$ and $\lim_{n \rightarrow \infty} \beta(n) = \bar{\beta}$. Then $\beta(n)$ satisfies

$$\beta_i(n) = \alpha_i(n), \quad \forall i \in L \setminus M \quad (12)$$

$$\min_{i \in I_{\text{up}}(\beta(n)) \cap M} F_i(\beta(n)) \geq \max_{i \in I_{\text{low}}(\beta(n)) \cap M} F_i(\beta(n)) \quad (13)$$

for all n . It is obvious from (12) that

$$\bar{\beta}_i = \bar{\alpha}_i, \quad \forall i \in L \setminus M.$$

Also, by applying the argument used in the proof of Proposition 2 to (13), we have

$$\min_{i \in I_{\text{up}}(\bar{\beta}) \cap M} F_i(\bar{\beta}) \geq \max_{i \in I_{\text{low}}(\bar{\beta}) \cap M} F_i(\bar{\beta}).$$

Therefore $\bar{\beta}$ belongs to $\Gamma_M(\bar{\alpha})$ which implies that $\Gamma_M(\alpha)$ is closed at $\bar{\alpha}$. Since $\bar{\alpha}$ can be any point in S , we can conclude that $\Gamma_M(\alpha)$ is closed on S . ■

Lemma 8: The point-to-set map $\mathbf{A}(\alpha)$ defined by (11) is closed on $S \setminus \Omega^{(\tau, \epsilon)}$.

Proof: Let $\{\alpha(n)\}_{n=0}^\infty$ be any sequence such that $\alpha(n) \in S \setminus \Omega^{(\tau, \epsilon)}$, $\forall n$ and $\lim_{n \rightarrow \infty} \alpha(n) = \bar{\alpha} \in S \setminus \Omega^{(\tau, \epsilon)}$. Let $\{\beta(n)\}_{n=0}^\infty$ be any sequence such that $\beta(n) \in \mathbf{A}(\alpha(n))$, $\forall n$ and $\lim_{n \rightarrow \infty} \beta(n) = \bar{\beta}$. Then it follows from Lemma 6 that there exists a positive integer n_1 such that

$$\beta(n) \in \cup_{M \in V_q(\bar{\alpha})} \Gamma_M(\alpha(n)), \quad \forall n \geq n_1. \quad (14)$$

As shown in Lemma 7, the point-to-set map $\Gamma_M(\alpha)$ is closed at $\bar{\alpha}$ for each $M \in V_q(\bar{\alpha})$. Moreover, we easily see from Lemma 1 that $\cup_{M \in V_q(\bar{\alpha})} \Gamma_M(\alpha)$ is closed at $\bar{\alpha}$. This result, together with (14), indicates that $\bar{\beta} \in \cup_{M \in V_q(\bar{\alpha})} \Gamma_M(\bar{\alpha})$. Therefore $\mathbf{A}(\alpha)$ is closed at $\bar{\alpha} \in S \setminus \Omega^{(\tau, \epsilon)}$. Since $\bar{\alpha}$ can be

any point in $S \setminus \Omega^{(\tau, \epsilon)}$, we can conclude that the point-to-set map $\mathbf{A}(\alpha)$ is closed on $S \setminus \Omega^{(\tau, \epsilon)}$. ■

Lemma 9: The objective function $W(\alpha)$ of Problem 2 is a descent function for the set of (τ, ϵ) -optimal solutions $\Omega^{(\tau, \epsilon)}$ and the point-to-set map $\mathbf{A}(\alpha)$ defined by (11).

Proof: Let β be any point belonging to $\mathbf{A}(\alpha)$. If $\alpha \notin \Omega^{(\tau, \epsilon)}$, there exists an $M \in V_q(\alpha)$ such that $\alpha \notin \Gamma_M(\alpha)$ and $\beta \in \Gamma_M(\alpha)$. This implies that $W(\beta) < W(\alpha)$. Therefore $W(\alpha)$ is a descent function for $\Omega^{(\tau, \epsilon)}$ and $\mathbf{A}(\alpha)$. ■

Now we are ready for giving the global convergence theorem for Algorithm 1, which is the main result of this paper.

Theorem 2: Let $\{\alpha(k)\}_{k=0}^\infty$ be the sequence generated by Algorithm 1. If the working set $L_B(k)$ contains at least one (τ, ϵ) -violating pair at $\alpha(k)$ for all k , then any convergent subsequence of $\{\alpha(k)\}_{k=0}^\infty$ has a limit in $\Omega^{(\tau, \epsilon)}$.

Proof: From the definition of the point-to-set map $\mathbf{A}(\alpha)$ in (11) and the assumption on the working set $L_B(k)$, it is apparent that $\alpha(k+1) \in \mathbf{A}(\alpha(k))$ for all k . The sequence $\{\alpha(k)\}_{k=0}^\infty$ belongs to S which is compact. As shown in Lemma 8, the map $\mathbf{A}(\alpha)$ is closed on $S \setminus \Omega^{(\tau, \epsilon)}$. Also, as shown in Lemma 9, the objective function $W(\alpha)$ is a descent function for $\Omega^{(\tau, \epsilon)}$ and $\mathbf{A}(\alpha)$. Therefore, we can conclude from Theorem 1 that any convergent subsequence of $\{\alpha(k)\}_{k=0}^\infty$ has a limit in $\Omega^{(\tau, \epsilon)}$. ■

The following theorem is immediately derived from Theorem 2 and Proposition 3.

Theorem 3: If the working set $L_B(k)$ contains at least one (τ, ϵ) -violating pair at $\alpha(k)$ for all k , then Algorithm 1 stops at $\Omega^{(\tau, \epsilon)}$ within a finite number of iterations for any $\tau > 0$ and $\epsilon \in (0, C/2)$.

Remark 1: According to Proposition 4, we can make $\Omega^{(\tau, \epsilon)}$ as close as we want to Ω^* by setting τ and ϵ to sufficiently small positive numbers. In other words, we can make the limit of the sequence of solutions generated by Algorithm 1 as close as we want to an optimal solution of Problem 2.

Remark 2: In the above discussion, we have assumed that the subproblems can be solved exactly. However, this assumption is not necessarily required. In fact, if it is guaranteed that $W(\alpha(k+1))$ is less than $W(\alpha(k))$ as far as $L_B(k)$ contains at least one (τ, ϵ) -violating pair at $\alpha(k)$, then Theorems 2 and 3 still hold. This weaker condition will be useful in practical situations where the subproblems are solved numerically and thus only approximate solutions can be obtained.

Remark 3: If we employ (6) instead of (7) for the optimality condition, the global convergence of Algorithm 1 cannot be proved as above because the point-to-set map corresponding to (11) is not closed in this case.

Let us next consider the case where (τ, ϵ) -violating pairs are not always contained in $L_B(k)$. Theorem 2 does not hold in this case because $W(\alpha)$ is not a descent function for $\Omega^{(\tau, \epsilon)}$ and $\mathbf{A}(\alpha)$ defined by (11). However, if at least one (τ, ϵ) -violating pair is selected for the working set within a certain period of iterations, Algorithm 1 still has the convergence property. This is formally stated as follows.

Theorem 4: Let $\{\alpha(k)\}_{k=0}^\infty$ be the sequence generated by Algorithm 1. If there exists a positive integer m such that one of m sets $L_B(k), L_B(k+1), \dots, L_B(k+m-1)$ contains at least one (τ, ϵ) -violating pair at $\alpha(k)$ for all k , then

Algorithm 1 stops at $\Omega^{(\tau, \epsilon)}$ within a finite number of iterations for any $\tau > 0$ and $\epsilon \in (0, C/2)$.

Proof: Let us define the point-to-set map $\mathbf{A}^m(\alpha)$ as

$$\mathbf{A}^m(\alpha) = \begin{cases} \cup_{(M_1, \dots, M_m) \in V_q^m(\alpha)} \Gamma_{M_m} \cdots \Gamma_{M_1}(\alpha), & \text{if } \alpha \notin \Omega^{(\tau, \epsilon)} \\ \alpha, & \text{if } \alpha \in \Omega^{(\tau, \epsilon)} \end{cases}$$

where $V_q^m(\alpha)$ is the set of all sequences of the sets (M_1, M_2, \dots, M_m) such that 1) $M_i \subseteq L$, $i = 1, 2, \dots, m$, 2) $|M_i| \leq q$, $i = 1, 2, \dots, m$, and 3) at least one M_i contains at least one (τ, ϵ) -violating pair at α . Let $\{\tilde{\alpha}(k)\}_{k=0}^\infty$ be the sequence defined by $\tilde{\alpha}(k) = \alpha(mk)$. Then the sequence $\{\tilde{\alpha}(k)\}_{k=0}^\infty$ satisfies

$$\tilde{\alpha}(k+1) \in \mathbf{A}^m(\tilde{\alpha}(k)) \subseteq S, \quad \forall k.$$

By applying Lemmas 1, 2, 6 and 7 we can easily show that $\mathbf{A}^m(\alpha)$ is closed on $S \setminus \Omega^{(\tau, \epsilon)}$. Moreover, it follows from Property 3) of $V_q^m(\alpha)$ mentioned above that the objective function $W(\alpha)$ is a descent function for $\Omega^{(\tau, \epsilon)}$ and $\mathbf{A}^m(\alpha)$. Thus, by Theorem 1, any convergent subsequence of $\{\tilde{\alpha}(k)\}_{k=0}^\infty$ has a limit in $\Omega^{(\tau, \epsilon)}$. Since $S \setminus \Omega^{(\tau, \epsilon)}$ is closed for any $\tau > 0$ and $\epsilon \in (0, C/2)$, the convergent subsequence enters $\Omega^{(\tau, \epsilon)}$ within a finite number of iterations. ■

V. APPLICATION TO EXISTING LEARNING ALGORITHMS

In this section, we will discuss the global convergence of some existing decomposition methods by applying the results obtained in the previous section.

A. Generalized SMO Algorithm

Generalized SMO algorithm [9] is a special type of decomposition methods in which the optimality condition is given by (6) and the working set is composed of a pair of indices (i, j) violating (6) for $\alpha = \alpha(k)$, that is,

$$i \in I_{\text{up}}(\alpha(k)), j \in I_{\text{low}}(\alpha(k)), F_i(\alpha(k)) < F_j(\alpha(k)) - \tau.$$

It has already been proved that this algorithm always stops within a finite number of iterations [9], [11]. We will now show that almost the same result can be derived by using a theorem in the previous section. To do so, we need to modify the generalized SMO algorithm slightly; the optimality condition is given by (7) instead of (6), and the working set is composed of a (τ, ϵ) -violating pair at $\alpha(k)$ where τ and ϵ are sufficiently small positive numbers. Then it is easily seen from Theorem 3 that the algorithm stops within a finite number of iterations.

B. SVM^{light}

SVM^{light} proposed by Joachims [6] is one of the most widely used decomposition methods for SVMs. The optimality condition used in SVM^{light} is given by (6) and the working set selection is done in a systematic way as follows:

Algorithm 2: Given an even number $q(\leq l)$ and the current solution $\alpha(k) \in S$, execute the following procedures.

- 1) Sort $\{F_i(\alpha(k))\}_{i=1}^l$ in decreasing order. Let the list of subscripts of the sorted list be i_1, i_2, \dots, i_l .

- 2) Set $L_B(k) = \emptyset$, $v = 0$, $m = 1$ and $n = l$.
- 3) While $i_m \notin I_{\text{low}}(\alpha(k))$ and $m \leq l$, add 1 to m .
- 4) While $i_n \notin I_{\text{up}}(\alpha(k))$ and $n \geq 1$, subtract 1 from n .
- 5) If $m \geq n$ then stop. Otherwise, add 2 to v and add $\{m, n\}$ to the set $L_B(k)$.
- 6) If $v = q$ then stop. Otherwise, go to Step 3).

For more details on the working set selection of SVM^{light}, please consult [14].

As well as in case of the generalized SMO algorithm, we need to modify both the optimality condition and the working set selection in SVM^{light} in order to apply the results in the previous section; the optimality condition is given by (7) instead of (6) and the working set selection is done by the following algorithm:

Algorithm 3: Given an even number $q(\leq l)$ and the current solution $\alpha(k) \in S$, execute the following procedures.

- 1) Sort $\{F_i(\alpha(k))\}_{i=1}^l$ in decreasing order. Let the the list of subscripts of the sorted list be i_1, i_2, \dots, i_l .
- 2) Set $L_B(k) = \emptyset$, $v = 0$, $m = 1$ and $n = l$.
- 3) While $i_m \notin I_{\text{low}}^\epsilon(\alpha(k))$ and $m \leq l$, add 1 to m .
- 4) While $i_n \notin I_{\text{up}}^\epsilon(\alpha(k))$ and $n \geq 1$, subtract 1 from n .
- 5) If $m \geq n$ then stop. Otherwise, add 2 to v and add $\{m, n\}$ to the set $L_B(k)$.
- 6) If $v = q$ then stop. Otherwise, go to Step 3).

We will show that $L_B(k)$ obtained by Algorithm 3 contains at least one (τ, ϵ) -violating pair if $\alpha(k) \in S \setminus \Omega^{(\tau, \epsilon)}$. Recall that $\alpha \in \Omega^{(\tau, \epsilon)}$ if and only if (7) holds. Thus if $\alpha(k) \in S \setminus \Omega^{(\tau, \epsilon)}$, the following inequality holds.

$$\min_{i \in I_{\text{up}}^\epsilon(\alpha(k))} F_i(\alpha(k)) + \tau \leq \max_{i \in I_{\text{low}}^\epsilon(\alpha(k))} F_i(\alpha(k)) \quad (15)$$

Let m_1 be the value of m at the instance when Algorithm 3 first exits Step 3). Then $i_{m_1} \in I_{\text{low}}^\epsilon(\alpha(k))$ and $F_{i_{m_1}}(\alpha(k)) = \max_{i \in I_{\text{low}}^\epsilon(\alpha(k))} F_i(\alpha(k))$. Also, let n_1 be the value of n at the instance when Algorithm 3 first exits Step 4). Then $i_{n_1} \in I_{\text{up}}^\epsilon(\alpha(k))$ and $F_{i_{n_1}}(\alpha(k)) = \min_{i \in I_{\text{up}}^\epsilon(\alpha(k))} F_i(\alpha(k))$. It follows from (15) that $F_{n_1}(\alpha(k)) + \tau \leq F_{m_1}(\alpha(k))$ which implies that (n_1, m_1) is a (τ, ϵ) -violating pair at $\alpha(k)$ and that $m_1 < n_1$. Therefore $L_B(k)$ obtained by Algorithm 3 contains at least one (τ, ϵ) -violating pair at $\alpha(k)$.

By applying Theorem 3, we can conclude that SVM^{light} with the optimality condition (7) and the working set selection described by Algorithm 3 stops within a finite number of iterations after finding a (τ, ϵ) -optimal solution.

Remark 4: In the software package SVM^{light} the optimality condition and the working set selection are not implemented exactly as described in (6) and Algorithm 2, respectively, for practical reason. Difference is in the definition of $I_{\text{up}}(\alpha)$ and $I_{\text{low}}(\alpha)$. In the software, the following definitions are used:

$$I_{\text{up}}(\alpha) = \{i \mid \alpha_i < C - \epsilon, d_i = 1\} \cup \{i \mid \alpha_i > \epsilon, d_i = -1\},$$

$$I_{\text{low}}(\alpha) = \{i \mid \alpha_i < C - \epsilon, d_i = -1\} \cup \{i \mid \alpha_i > \epsilon, d_i = 1\}$$

where ϵ is a sufficiently small positive number. It is worth noting that these definitions are same as $I_{\text{up}}^\epsilon(\alpha)$ and $I_{\text{low}}^\epsilon(\alpha)$ except for the equal sign. The results of this paper show that the relaxation of the optimality condition has not only practical but also theoretical significance.

VI. CONCLUDING REMARKS

Global convergence property of decomposition methods for SVMs are studied. We have first introduced a relaxed optimality condition, and then proved a decomposition method stops within a finite number of iterations after finding an optimal solution if the working set selection satisfies a certain condition. We have also shown that the generalized SMO algorithm and SVM^{light} satisfy this condition and thus have the global convergence property. Since our new convergence theorems require little restriction on the working set selection method, the authors believe that they can be applied to a wide class of decomposition methods.

REFERENCES

- [1] V. Vapnik, *Statistical Learning Theory*. New York: Wiley, 1998.
- [2] B. Schölkopf, C. J. C. Burges, and A. J. Smola, Eds., *Advances in Kernel Methods: Support Vector Learning*. Cambridge, Massachusetts: MIT Press, 1999.
- [3] N. Cristianini and J. Shawe-Taylor, *An Introduction to Support Vector Machines and Other Kernel-Based Learning Methods*. Cambridge, UK: Cambridge University Press, 2000.
- [4] S. Haykin, *Neural Networks*. Upper Saddle River, NJ: Prentice Hall, 1999.
- [5] J. C. Platt, "Fast training of support vector machines using sequential minimal optimization," in *Advances in Kernel Methods: Support Vector Machines*, B. Schölkopf, C. Burges, and A. Smola, Eds. Cambridge, MA: MIT Press, 1998.
- [6] T. Joachims, "Making large-scale support vector machine learning practical," in *Advances in Kernel Methods: Support Vector Machines*, B. Schölkopf, C. Burges, and A. Smola, Eds. Cambridge, MA: MIT Press, 1998.
- [7] S. S. Keerthi, S. K. Shevade, C. Bhattacharyya, and K. R. K. Murthy, "Improvements to Platt's SMO algorithm for SVM classifier design," *Neural Computation*, vol. 13, pp. 637–649, 2001.
- [8] C.-W. Hsu and C.-J. Lin, "A simple decomposition method for support vector machines," *Machine Learning*, vol. 46, pp. 291–314, 2002.
- [9] S. S. Keerthi and E. G. Gilbert, "Convergence of a generalized SMO algorithm for SVM classifier design," *Machine Learning*, vol. 46, pp. 351–360, 2002.
- [10] C.-J. Lin, "Asymptotic convergence of an SMO algorithm without any assumption," *IEEE Trans. Neural Networks*, vol. 13, pp. 248–250, Jan. 2002.
- [11] N. Takahashi and T. Nishi, "Rigorous proof of termination of SMO algorithm for support vector machines," *IEEE Trans. Neural Networks*, vol. 16, no. 3, pp. 774–776, May 2005.
- [12] P.-H. Chen, R.-E. Fan, and C.-J. Lin, "A study on SMO-type decomposition methods for support vector machines," to appear in *IEEE Trans. Neural Networks*.
- [13] C.-C. Chang, C.-W. Hsu, and C.-J. Lin, "The analysis of decomposition methods for support vector machines," *IEEE Trans. Neural Networks*, vol. 11, no. 4, pp. 1003–1008, 2000.
- [14] C.-J. Lin, "On the convergence of the decomposition method for support vector machines," *IEEE Trans. Neural Networks*, vol. 12, pp. 1288–1298, Nov. 2001.
- [15] —, "A formal analysis of stopping criteria of decomposition methods for support vector machines," *IEEE Trans. Neural Networks*, vol. 13, pp. 1045–1052, Sept. 2002.
- [16] N. List and H. U. Simon, "A general convergence theorem for the decomposition method," in *Proceedings of the 17th Annual Conference on Learning Theory*, 2004, pp. 363–377.
- [17] D. G. Luenberger, *Linear and Nonlinear Programming*. Reading, MA: Addison-Wesley Publishing, 1989.