

A Formal Analysis of Stopping Criteria of Decomposition Methods for Support Vector Machines

Chih-Jen Lin, *Member, IEEE*

Abstract—In a previous paper, we proved the convergence of a commonly used decomposition method for support vector machines (SVMs). However, there is no theoretical justification about its stopping criterion, which is based on the gap of the violation of the optimality condition. It is essential to have the gap asymptotically approach zero, so we are sure that existing implementations stop in a finite number of iterations after reaching a specified tolerance. Here, we prove this result and illustrate it by two extensions: ν -SVM and a multiclass SVM by Crammer and Singer. A further result shows that, in final iterations of the decomposition method, only a particular set of variables are still being modified. This supports the use of the shrinking and caching techniques in some existing implementations. Finally, we prove the asymptotic convergence of a decomposition method for this multiclass SVM. Discussions on the difference between this convergence proof and the one in another paper by Lin are also included.

Index Terms—Asymptotic convergence, decomposition methods, stopping criteria, support vector machines (SVMs).

I. INTRODUCTION

GIVEN a training set of instance-label pairs $(x_i, y_i), i = 1, \dots, l$ where $x_i \in R^n$ and $y \in \{1, -1\}^l$, the support vector machines (SVMs) [3], [13] require the solution of the following optimization problem:

$$\begin{aligned} \min_{w, b, \xi} \quad & \frac{1}{2} w^T w + C \sum_{i=1}^l \xi_i \\ \text{subject to} \quad & y_i (w^T \phi(x_i) + b) \geq 1 - \xi_i \\ & \xi_i \geq 0, \quad i = 1, \dots, l. \end{aligned} \quad (1)$$

Here, training vectors x_i are mapped into a higher (maybe infinite) dimensional space by the function ϕ . Then SVM finds a linear separating hyperplane with the maximal margin in this higher dimensional space. $C > 0$ is the penalty parameter of the error term. As the number of variables becomes large after mapping the data, practically we solve the dual problem

$$\begin{aligned} \min_{\alpha} \quad & f(\alpha) = \frac{1}{2} \alpha^T Q \alpha - e^T \alpha \\ \text{subject to} \quad & 0 \leq \alpha_i \leq C, \quad i = 1, \dots, l \\ & y^T \alpha = 0 \end{aligned} \quad (2)$$

where e is the vector of all ones, C becomes the upper bound of all variables α , and Q is an l by l positive semidefinite matrix. Note that $Q_{ij} \equiv y_i y_j K(x_i, x_j)$ where $K(x_i, x_j) \equiv \phi(x_i)^T \phi(x_j)$ is called the kernel function. Then $w = \sum_{i=1}^l \alpha_i y_i \phi(x_i)$ and

$$\text{sgn}(w^T \phi(x) + b) = \text{sgn} \left(\sum_{i=1}^l \alpha_i y_i K(x_i, x) + b \right)$$

is the decision function.

Due to the density of the matrix Q , currently the decomposition method is one of the major methods to solve SVM (e.g., [6], [10], and [11]). It is an iterative process, wherein each iteration the index set of variables is separated to two sets B and N , where B is the working set. Then, in that iteration, variables corresponding to N are fixed, while a subproblem on variables corresponding to B is minimized.

Practically, we need a stopping condition for the decomposition method. Such a criterion usually uses the information of the Karush–Kuhn–Tucker (KKT) condition, that is, the optimality condition of (2): If α is an optimal solution of (2), there is a number b and two nonnegative vectors λ and μ such that

$$\begin{aligned} \nabla f(\alpha) + by = \lambda - \mu, \quad & \lambda_i \alpha_i = 0, \quad \mu_i (C - \alpha)_i = 0 \\ & \lambda_i \geq 0, \quad \mu_i \geq 0, \quad i = 1, \dots, l. \end{aligned}$$

This is usually rewritten as

$$\nabla f(\alpha)_i + by_i \geq 0, \quad \text{if } \alpha_i = 0 \quad (3a)$$

$$\nabla f(\alpha)_i + by_i \leq 0, \quad \text{if } \alpha_i = C \quad (3b)$$

$$\nabla f(\alpha)_i + by_i = 0, \quad \text{if } 0 < \alpha_i < C. \quad (3c)$$

As $C > 0$, we can further reformulate it as

$$\nabla f(\alpha)_i + by_i \geq 0, \quad \text{if } \alpha_i < C \quad (4a)$$

$$\nabla f(\alpha)_i + by_i \leq 0, \quad \text{if } \alpha_i > 0. \quad (4b)$$

Since $y_i = \pm 1$, by expressing inequalities of (4) as lower and upper bounds of b , this KKT condition is equivalent to

$$\begin{aligned} m(\alpha) &= \max \left(\max_{\alpha_i < C, y_i = 1} -\nabla f(\alpha)_i, \max_{\alpha_i > 0, y_i = -1} \nabla f(\alpha)_i \right) \\ &\leq \min \left(\min_{\alpha_i < C, y_i = -1} \nabla f(\alpha)_i, \min_{\alpha_i > 0, y_i = 1} -\nabla f(\alpha)_i \right) \\ &= M(\alpha). \end{aligned} \quad (5)$$

Manuscript received June 1, 2001; revised January 16, 2002. This work was supported in part by the National Science Council of Taiwan under Grant NSC 90-2213-E-002-111.

The author is with the Department of Computer Science and Information Engineering, National Taiwan University, Taipei 106, Taiwan (e-mail: cjlin@csie.ntu.edu.tw).

Publisher Item Identifier S 1045-9227(02)04437-5.

Let α^k be the solution at the k th iteration. If α^k is not an optimal solution then $m(\alpha^k) > M(\alpha^k)$. Hence, a natural stopping criterion might be

$$m(\alpha^k) \leq M(\alpha^k) + \epsilon \quad (6)$$

where ϵ is a stopping tolerance. We can see that (5) is a much simpler way of describing the KKT condition. Some existing working set selections also follow from identifying elements violating (5). More importantly, unlike earlier approaches where b is calculated using (3c), the condition on free variables, we do not have to worry if there are free variables in the final solution or not. If all variables are at bounds, using (5), b can be simply calculated as $(M(\alpha) + m(\alpha))/2$. Such a stopping criterion has been derived and used in, for example, [1] and [7].

In an earlier work [8] on the convergence of the decomposition method proposed in the software SVM^{light} [6], we focused on proving that any limit point of $\{\alpha^k\}$ is an optimal solution of (2). However, such results do not directly support the validity of using (6) as the stopping criterion. To be more precise, even though if we have $\lim_{k \rightarrow \infty} \alpha^k = \bar{\alpha}$ which is an optimal solution, directly from the definition of $m(\alpha^k)$ and $m(\bar{\alpha})$, we may not have $\lim_{k \rightarrow \infty} m(\alpha^k) = m(\bar{\alpha})$. A similar problem happens for $M(\alpha^k)$ and $M(\bar{\alpha})$. A reason is that it may be possible that $\alpha_i^k > 0, \forall k$ but $\lim_{k \rightarrow \infty} \alpha_i^k = \bar{\alpha}_i = 0$. Therefore, we worry about the situation that $\{\alpha^k\}$ converges to an optimal solution but $\lim_{k \rightarrow \infty} m(\alpha^k) - M(\alpha^k) > 0$. Then, the decomposition implementation never stops by using (6) as the criterion. In Section II, we prove that this situation will never happen.

Note that if the size of the working set is restricted to two, the finite termination of using the stopping criterion (6) has been proved in [7]. To be more precise, they prove that for any tolerance ϵ , the algorithm stops in a finite number of iterations. However, as discussed previously, their result does not imply the asymptotic convergence (i.e., any limit point of $\{\alpha^k\}$ is an optimum). On the other hand, for more general analyses on the stopping criterion (6), we will use the asymptotic convergence which has been proved in [8].

In Section II, we also prove that most bounded variables are identified after finite steps so final iterations focus on a particular set of variables. This analysis supports the use of shrinking and caching techniques in the decomposition method. Section III then shows some extensions on a more complicated optimization formulation. We use two examples to illustrate our results: ν -SVM [12] and a multiclass SVM in [4]. We then also prove the asymptotic convergence of a decomposition method for this multiclass SVM. Discussions on the difference between this convergence proof and the one in [8] are also included.

II. MAIN RESULTS ON STOPPING CRITERIA

Here, we consider a more general problem

$$\begin{aligned} & \min_{\alpha} f(\alpha) \\ & \text{subject to } o_i \leq \alpha_i \leq u_i, \quad i = 1, \dots, l \\ & \quad \quad \quad y^T \alpha = 0 \end{aligned} \quad (7)$$

where $f(\alpha)$ is any continuously convex differentiable function, $y_i = \pm 1, i = 1, \dots, l$, and o and u are lower and upper bounds,

respectively. Following the requirement that $C > 0$, here we consider only the following situation:

$$-\infty < o_i < u_i < \infty. \quad (8)$$

If $o_i = u_i$ for some i , $\alpha_i = o_i = u_i$ so we can remove such variables before solving the problem by decomposition methods. We then define generalized $m(\alpha^k)$ and $M(\alpha^k)$

$$m(\alpha^k) \equiv \max \left(\max_{\alpha_i^k < u_i, y_i = 1} -\nabla f(\alpha^k)_t, \max_{\alpha_i^k > o_i, y_i = -1} \nabla f(\alpha^k)_t \right) \quad (9)$$

and

$$M(\alpha^k) \equiv \min \left(\min_{\alpha_i^k < u_i, y_i = -1} \nabla f(\alpha^k)_t, \min_{\alpha_i^k > o_i, y_i = 1} -\nabla f(\alpha^k)_t \right). \quad (10)$$

We denote $\arg m(\alpha^k)$ the set of indexes whose $-y_i \nabla f(\alpha^k)_i$ are the same as $m(\alpha^k)$. A similar definition goes for $\arg M(\alpha^k)$. Thus, if α^k is not an optimal solution yet, $m(\alpha^k) > M(\alpha^k)$ so $\arg m(\alpha^k) \cap \arg M(\alpha^k) = \emptyset$.

The following theorem shows the validity of the stopping criterion (6).

Theorem II.1: Assume a decomposition method for solving (7) satisfies the following conditions.

- 1) $m(\alpha^k) > M(\alpha^k), \forall k$.
- 2) At least one element of $\arg m(\alpha^k)$ and one element of $\arg M(\alpha^k)$ are included in the working set of each iteration.
- 3) For variables considered in (10), if (α_i^k, α_j^k) are any two of them such that

a) α_i^k satisfies

$$\alpha_i^k = o_i, \quad y_i = -1 \quad \text{or} \quad \alpha_i^k = u_i, \quad y_i = 1;$$

b) α_j^k satisfies

$$\alpha_j^k < u_j, \quad y_j = -1 \quad \text{or} \quad \alpha_j^k > o_j, \quad y_j = 1;$$

c)

$$-y_i \nabla f(\alpha^k)_i > -y_j \nabla f(\alpha^k)_j;$$

then, if α_i^k is in the working set, α_j^k must be selected as well.

- 4) Similar to Condition 3), we assume that analogous conditions hold for variables considered in (9).
- 5) $\{\alpha^k\}$ converges to an optimal solution $\bar{\alpha}$ of (7).

Then

$$\lim_{k \rightarrow \infty} (m(\alpha^k) - M(\alpha^k)) = 0. \quad (11)$$

Proof: We prove the theorem by contradiction and, thus, let us assume that the result (11) is wrong. Then, with Condition 1), there is an infinite set $\bar{\mathcal{K}}$ and a $\Delta > 0$ such that

$$m(\alpha^k) \geq M(\alpha^k) + \Delta \quad \forall k \in \bar{\mathcal{K}}. \quad (12)$$

As $\bar{\mathcal{K}}$ has infinitely many elements but the number of pairs of variables is finite, from Condition 2), there are indexes i and j and an infinite subsequence such that $i \in \arg m(\alpha^k)$ and

$j \in \arg M(\alpha^k)$ are both in the working set. Without loss of generality, we can consider only the case that there is $\hat{\mathcal{K}} \subset \bar{\mathcal{K}}$ such that $y_i = y_j = -1$ and for all $k \in \hat{\mathcal{K}}$

$$m(\alpha^k) = \nabla f(\alpha^k)_i \quad \text{and} \quad M(\alpha^k) = \nabla f(\alpha^k)_j. \quad (13)$$

Thus, from the definition of $m(\alpha^k)$ and $M(\alpha^k)$ in (9) and (10)

$$\alpha_i^k > o_i \quad \text{and} \quad \alpha_j^k < u_j \quad \forall k \in \hat{\mathcal{K}}. \quad (14)$$

Since f is continuously differentiable

$$\lim_{k \rightarrow \infty} \nabla f(\alpha^k)_i = \nabla f(\bar{\alpha})_i \quad \text{and} \quad \lim_{k \rightarrow \infty} \nabla f(\alpha^k)_j = \nabla f(\bar{\alpha})_j. \quad (15)$$

From (12), (13), and (15), we have

$$\nabla f(\bar{\alpha})_i > \nabla f(\bar{\alpha})_j. \quad (16)$$

If we have an infinite subset of $\hat{\mathcal{K}}$ such that $\alpha_i^{k+1} > o_i$ and $\alpha_j^{k+1} < u_j$, then the KKT condition of the subproblem implies

$$\nabla f(\alpha^{k+1})_j \geq \nabla f(\alpha^{k+1})_i.$$

Since $\{\alpha^k\}$ is a convergent sequence, taking the limit we have

$$\nabla f(\bar{\alpha})_i \leq \nabla f(\bar{\alpha})_j \quad (17)$$

which contradicts (16).

Therefore, we have that

$$\alpha_i^{k+1} = o_i \quad \text{or} \quad \alpha_j^{k+1} = u_j, \quad \text{after } k \in \hat{\mathcal{K}} \text{ is large enough.} \quad (18)$$

Because of (18) and (14), there is an infinite set \mathcal{L} such that for all $k \in \mathcal{L}$

$$\begin{aligned} \alpha_i^k = o_i, \quad \alpha_j^k < u_j, \quad \text{or} \quad \alpha_i^k > o_i, \quad \alpha_j^k = u_j, \quad \text{or} \\ \alpha_i^k = o_i, \quad \alpha_j^k = u_j \end{aligned} \quad (19)$$

and

$$\alpha_i^{k+1} > o_i, \quad \alpha_j^{k+1} < u_j. \quad (20)$$

For the first case of (19), α_i^k is selected in the working set and then modified. However, since (16) and Condition 5) imply

$$\nabla f(\alpha^k)_i > \nabla f(\alpha^k)_j$$

after k is large enough, with $\alpha_j^k < u_j$ in (19), and Condition 3) of this theorem, α_j^k is also in the working set of the k th iteration where $k \in \mathcal{L}$. With (20), from the KKT condition of the subproblem

$$\nabla f(\alpha^{k+1})_i \leq \nabla f(\alpha^{k+1})_j. \quad (21)$$

The situation for the second case is similar. For the third case, both α_i^k and α_j^k are modified so i and j are in the working set. Hence, (21) is also valid.

Therefore, we have (21) for all $k \in \mathcal{L}$. As k goes to infinity, we again obtain (17) which contradicts (16). Thus, the assumption (12) is wrong so the proof is complete. ■

Note that Conditions 2)-4) of Theorem II.1 are requirements on the working set selection. We list conditions instead of focusing on a particular working set selection so that more flexible selections may be used.

We now check that the working set selection of SVM^{light} satisfies the Conditions 2)-4) of Theorem II.1. If q , an even number, is the size of the working set, $q/2$ indexes are sequentially selected from elements that satisfy $\alpha_t^k < u_i, y_t = 1$ or $\alpha_t^k > o_i, y_t = -1$ so that

$$\begin{aligned} -y_{i_1} \nabla f(\alpha^k)_{i_1} &\geq -y_{i_2} \nabla f(\alpha^k)_{i_2} \geq \dots \\ &\geq -y_{i_{q/2}} \nabla f(\alpha^k)_{i_{q/2}} \quad \text{and } i_1 \in \arg m(\alpha^k). \end{aligned} \quad (22)$$

The other $q/2$ indexes are sequentially selected from elements which satisfy $\alpha_t^k < u_i, y_t = -1$ or $\alpha_t^k > o_i, y_t = 1$ such that

$$\begin{aligned} -y_{j_{q/2}} \nabla f(\alpha^k)_{j_{q/2}} &\geq \dots \geq -y_{j_1} \nabla f(\alpha^k)_{j_1} \quad \text{and} \\ j_1 &\in \arg M(\alpha^k). \end{aligned} \quad (23)$$

It can be clearly seen that directly from (22) and (23) these conditions are satisfied.

An interesting note is that this working set selection was originally derived from the concept of feasible directions in constrained optimization but not from the violation of the KKT condition.

Regarding the global convergence of $\{\alpha^k\}$ which is the Condition 5), unfortunately, we prove only a weaker result in [8]: under a minor assumption every limit point of convergent subsequences is an optimal solution. However, results in [8] do imply the global convergence if (2) has a unique optimal solution. Then Theorem II.1 can be applied. For example, if Q is positive definite, the solution of (2) is unique.

In the following we will show that for the algorithm used by SVM^{light}, Theorem II.1 is valid without needing the global convergence of $\{\alpha^k\}$. However, we still need the property proved in [8] that any limit point is an optimum. As all other conditions of Theorem II.1 are satisfied for this particular working set selection, the only remaining assumption is a minor one used in [8].

Theorem II.2: Under [8, Assumption IV.1], the decomposition method using (22) and (23) for selecting the working set has that if $m(\alpha^k) - M(\alpha^k) > 0, \forall k$, then

$$\lim_{k \rightarrow \infty} (m(\alpha^k) - M(\alpha^k)) = 0. \quad (24)$$

Proof: We also prove the theorem by contradiction. However, in addition to assuming an infinite sequence such that (12) holds, using results in [8], we further consider one of its convergent subsequence whose limit $\bar{\alpha}$ is an optimum of (2). Note that here any infinite sequence in the feasible region of (2) has at least one convergent subsequence because (8) implies that the feasible region is compact. Then, with [8, Assumption IV.1], the limit point is an optimum. Therefore, we can consider an infinite set $\bar{\mathcal{K}}$ and $\Delta > 0$ such that (12) holds and $\lim_{k \in \bar{\mathcal{K}}} \alpha^k = \bar{\alpha}$ is an optimum of (2).

Then, until (16), the proof is similar to that of Theorem II.1. Of course, $\lim_{k \rightarrow \infty}$ in (15) must be replaced by $\lim_{k \rightarrow \infty, k \in \bar{\mathcal{K}}}$

Remember that we consider the case of $y_i = y_j = -1$. We then claim that for all $k \in \hat{\mathcal{K}}$ large enough

$$\alpha_i^{k-1} > o_i \quad \text{and} \quad \alpha_j^{k-1} < u_j. \quad (25)$$

If (25) is wrong

$$\begin{aligned} \alpha_i^{k-1} = o_i, \quad \alpha_j^{k-1} < u_j, \quad \text{or} \quad \alpha_i^{k-1} > o_i \\ \alpha_j^{k-1} = u_j, \quad \text{or} \quad \alpha_i^{k-1} = o_i, \quad \alpha_j^{k-1} = u_j. \end{aligned} \quad (26)$$

For the first case, α_i^{k-1} is selected in the working set and modified. Since $\{\alpha^k\}, k \in \bar{\mathcal{K}}$ is a convergent subsequence, [8, Th. IV.3] implies that $\{\alpha^{k-1}\}, k \in \bar{\mathcal{K}}$ also converges to $\bar{\alpha}$. Thus, after $k \in \bar{\mathcal{K}}$ is large enough, (16) implies

$$-y_i \nabla f(\alpha^{k-1})_i > -y_j \nabla f(\alpha^{k-1})_j$$

Hence, (22) and (23) imply that α_j^{k-1} is selected in the working set as well. Therefore, from the KKT condition of the subproblem at the $(k-1)$ st iteration

$$\nabla f(\alpha^k)_i \leq \nabla f(\alpha^k)_j$$

which is impossible after $k \in \bar{\mathcal{K}}$ is large enough. The situation for other cases of (26) is similar. Therefore, (25) is correct.

Since [8, Th. IV.3] shows that for any given s , $\{\alpha^{k-s}\}, k \in \bar{\mathcal{K}}$ converges to the same point $\bar{\alpha}$ as $\{\alpha^k\}, k \in \bar{\mathcal{K}}$, using the same argument above we have that for any given s , after $k \in \bar{\mathcal{K}}$ is large enough

$$\begin{aligned} \alpha_i^{k-s} > o_i \quad \text{and} \quad \alpha_j^{k-s} < u_j, \dots, \alpha_i^k > o_i \quad \text{and} \\ \alpha_j^k < u_j, \quad -y_i \nabla f(\alpha^{k-s})_i > -y_j \nabla f(\alpha^{k-s})_j, \dots, \\ -y_i \nabla f(\alpha^k)_i > -y_j \nabla f(\alpha^k)_j. \end{aligned} \quad (27)$$

Consider $s = 2l$. Using the same counting procedure in [8, Th. IV.5], we can show that at some $k' \in \{k-2l, \dots, k-1\}$, $\alpha_i^{k'}$ and $\alpha_j^{k'}$ are both selected in the working set so the KKT condition of the subproblem again shows

$$f(\alpha^{k'+1})_i \leq f(\alpha^{k'+1})_j.$$

This is in contradiction to (27).

Therefore, the assumption (12) is wrong so

$$\lim_{k \rightarrow \infty} (m(\alpha^k) - M(\alpha^k)) = 0.$$

A special case of the working set selection using (22) and (23) is to restrict the size of the working set to two. That is, only one element of $\arg m(\alpha^k)$ and one element of $\arg M(\alpha^k)$ are included in the working set. This is an algorithm discussed in [7] and used in, for example, LIBSVM [1]. For this special case, [9] proves that [8, Assumption IV.1] is not necessary. Hence Theorem II.2 is valid without needing any assumption, exactly the same as results from [7].

Next, based on the above results, we show that after k is large enough, only elements whose $-y_i \nabla f(\bar{\alpha})_i$ are $m(\bar{\alpha})$ or $M(\bar{\alpha})$ can still be modified. For simplification, we only show results extended from Theorem II.1.

Theorem II.3: Under the same assumptions as Theorem II.1, we have the following result:

For any $\bar{\alpha}_i$ whose corresponding $-y_i \nabla f(\bar{\alpha})_i$ is neither $m(\bar{\alpha})$ nor $M(\bar{\alpha})$, after k is large enough, α_i^k is always at a bound which is equal to $\bar{\alpha}_i$.

Proof: First, we know that from the KKT condition, if $-y_i \nabla f(\bar{\alpha})_i$ is neither $m(\bar{\alpha})$ nor $M(\bar{\alpha})$, $\bar{\alpha}_i$ is at a bound. Without loss of generality, we consider

$$\bar{\alpha}_i = o_i \quad \text{with} \quad y_i = -1 \quad \text{and} \quad \nabla f(\bar{\alpha})_i > M(\bar{\alpha}). \quad (28)$$

If the result of this theorem is wrong, $\alpha_i^k > o_i$ happens infinitely many times. Therefore, from (28), there is an infinite set $\bar{\mathcal{K}}$ and a $\Delta > 0$ such that

$$m(\alpha^k) \geq \nabla f(\alpha^k)_i > M(\bar{\alpha}) + \Delta \quad \forall k \in \bar{\mathcal{K}}. \quad (29)$$

Now, for any $j \in \arg M(\bar{\alpha})$, we have $\bar{\alpha}_j < u_j, y_j = -1$ or $\bar{\alpha}_j > o_j, y_j = 1$. Therefore, $\alpha_j^k < u_j, y_j = -1$ or $\alpha_j^k > o_j, y_j = 1$ after k is large enough. Since $\{\alpha^k\}$ is a convergent sequence, $\{\alpha^k\}$ is in a compact region. With $M(\alpha^k) \leq -y_j \nabla f(\alpha^k)_j$, there is an infinite subset $\hat{\mathcal{K}}$ of $\bar{\mathcal{K}}$ such that $\lim_{k \in \hat{\mathcal{K}}} M(\alpha^k)$ exists and

$$\lim_{k \in \hat{\mathcal{K}}} M(\alpha^k) \leq M(\bar{\alpha}). \quad (30)$$

Hence, (29) and (30) imply

$$\lim_{k \in \hat{\mathcal{K}}} m(\alpha^k) - M(\alpha^k) \neq 0 \quad (31)$$

which contradicts Theorem II.1. This completes the proof. ■

Therefore, after k is large enough, only elements in

$$\{t \mid -y_t \nabla f(\bar{\alpha})_t = m(\bar{\alpha}) = M(\bar{\alpha})\} \quad (32)$$

can still be possibly modified.

This analysis supports the use of shrinking techniques [6] in the decomposition method as in final iterations it is possible that most variables are not changed any more. That is, after identifying some variables which may be at the bounds eventually, we temporarily remove them and solve a smaller optimization problem. Though a final check is still needed, in general, the training time can be largely saved.

Caching is another popular technique employed in implementations of decomposition methods. We store recently used kernel elements in the computer memory in order to save the number of kernel evaluations. Note that in each iteration as q variables are updated, basically q columns of the Hessian matrix have to be involved for updating the gradient $\nabla f(\alpha^k)$. Theorem II.3 supports this caching strategy as it shows that in final iterations only some particular columns of the kernel matrix are still needed.

For the algorithm used in SVM^{light}, Theorem II.3 becomes as follows:

Theorem II.4: Under [8, Assumption IV.1], if (22) and (23) are used for selecting the working set and $\bar{\alpha}$ is the limit point of any convergent subsequence $\{\alpha^k\}, k \in \mathcal{K}$, we have the following result.

For any $\bar{\alpha}_i$ whose corresponding $-y_i \nabla f(\bar{\alpha})_i$ is neither $m(\bar{\alpha})$ nor $M(\bar{\alpha})$, after $k \in \mathcal{K}$ is large enough, α_i^k is always at a bound which is equal to $\bar{\alpha}_i$.

The proof is nearly the same as that for Theorem II.3.

To conclude this section we note that (7) is a more general formulation so results in this section apply to different formulations discussed in [8] such as support vector regression and one-class SVM.

III. EXTENSIONS

In this section, we consider a more general problem

$$\begin{aligned} & \min_{\alpha} f(\alpha) \\ & \text{subject to} \quad \sum_{m=1}^{r_i} y_i^m \alpha_i^m = \Delta_i \\ & \quad \alpha_i^m \leq \alpha_i^m \leq u_i^m, \quad m = 1, \dots, r_i, \quad i = 1, \dots, l \end{aligned} \quad (33)$$

where $y_i^m = \pm 1$ and $-\infty < \alpha_i^m < u_i^m < \infty$. Therefore, there are $\sum_{i=1}^l r_i$ variables and l linear equality constraints. Hence, it is like there are l groups of variables where each one satisfies a linear constraint. The KKT condition requires that there are b_1, \dots, b_l such that for all $i = 1, \dots, l, m = 1, \dots, r_i$

$$\begin{aligned} \nabla f(\alpha)_i^m + b_i y_i^m & \geq 0 \quad \text{if } \alpha_i^m < u_i^m \\ & \leq 0 \quad \text{if } \alpha_i^m > o_i^m \end{aligned}$$

where $\nabla f(\alpha)_i^m$ means $\partial f(\alpha)/\partial \alpha_i^m$. We can rewrite the KKT condition as

$$\begin{aligned} & m(\alpha)_i \\ & = \max \left(\max_{\alpha_i^m < u_i^m, y_i^m = 1} -\nabla f(\alpha)_i^m, \max_{\alpha_i^m > o_i^m, y_i^m = -1} \nabla f(\alpha)_i^m \right) \\ & \leq \min \left(\min_{\alpha_i^m < u_i^m, y_i^m = -1} \nabla f(\alpha)_i^m, \min_{\alpha_i^m > o_i^m, y_i^m = 1} -\nabla f(\alpha)_i^m \right) \\ & = M(\alpha)_i, \quad i = 1, \dots, l. \end{aligned} \quad (34)$$

By defining

$$m(\alpha) \equiv \max_i m(\alpha)_i \quad \text{and} \quad M(\alpha) \equiv \min_i M(\alpha)_i$$

the stopping criterion can be

$$m(\alpha^k) - M(\alpha^k) \leq \epsilon \quad (35)$$

where ϵ is the stopping tolerance.

Similar to the situation in Section II, we can define two sets $\arg m(\alpha^k)$ and $\arg M(\alpha^k)$. With some minor modifications on Conditions 2)–4), we can have results similar to Theorem II.1.

Theorem III.1: Assume all conditions of Theorem II.1 hold with Conditions 2)–4) replaced by the following conditions.

- 2' In each iteration, if the i th group has the largest $m(\alpha^k)_i - M(\alpha^k)_i$, then at least one of $\arg m(\alpha^k)_i$ and one of $\arg M(\alpha^k)_i$ are included in the working set.
- 3' In each iteration, variables of the group with the largest $m(\alpha^k)_i - M(\alpha^k)_i$ satisfy Condition 3).
- 4' In each iteration, variables of the group with the smallest $m(\alpha^k)_i - M(\alpha^k)_i$ satisfy Condition 4).

Then

$$\lim_{k \rightarrow \infty} m(\alpha^k) - M(\alpha^k) = 0. \quad (36)$$

Some SVM formulations are of this form. We will give two examples: ν -SVM and a multiclass SVM by Crammer and Singer. The ν -SVM [12] can be written as the following form [2]:

$$\begin{aligned} & \min_{\alpha} \quad \frac{1}{2} \alpha^T Q \alpha \\ & \text{subject to} \quad \sum_{m=1}^{r_1} \alpha_1^m = \frac{\nu}{2}, \quad \sum_{m=1}^{r_2} \alpha_2^m = \frac{\nu}{2} \\ & \quad 0 \leq \alpha_1^m \leq \frac{1}{l}, \quad m = 1, \dots, r_1 \\ & \quad 0 \leq \alpha_2^m \leq \frac{1}{l}, \quad m = 1, \dots, r_2 \end{aligned}$$

where $\nu \in [0, 1]$ is a parameter to adjust the number of support vectors and training errors, r_1 and r_2 are numbers of training data in two classes, and $l = r_1 + r_2$. A stopping criterion like (35) has been used in the experiment of [2] which implemented a modified decomposition method from the one in [8]. We can easily check that Conditions 2'–4' of Theorem III.1 are satisfied. Regarding the convergence of $\{\alpha^k\}$, though we have not explicitly written down the proof, we conjecture that the same results in [8] that every limit point is an optimal solution should still apply.

Another example is a formulation for multiclass SVM by Crammer and Singer [4]

$$\begin{aligned} & \min_{\alpha} \quad \frac{1}{2} \alpha^T (K \otimes I) \alpha + \sum_{i=1}^l \bar{e}_i^T \alpha_i \\ & \text{subject to} \quad \sum_{m=1}^r \alpha_i^m = 0 \\ & \quad \alpha_i^m \leq C_{\bar{y}_i}^m, \quad i = 1, \dots, l, \quad m = 1, \dots, r \end{aligned} \quad (37a)$$

$$(37b)$$

where \otimes is the Kronecker product, r is the number of classes, K is an l by l kernel matrix, I is an r by r identity matrix, each \bar{e}_i is an r by 1 constant vector, and $\bar{y}_i \in \{1, \dots, r\}$ is the label of the i th data, and

$$C_{\bar{y}_i}^m = \begin{cases} 0, & \text{if } \bar{y}_i \neq m \\ C, & \text{if } \bar{y}_i = m. \end{cases}$$

Here, α is an rl by 1 vector variable and we denote

$$\alpha \equiv [\alpha_1^1, \dots, \alpha_1^r, \dots, \alpha_l^1, \dots, \alpha_l^r]^T$$

and

$$\alpha_i \equiv [\alpha_i^1, \dots, \alpha_i^r]^T, \quad i = 1, \dots, l.$$

Hence, the stopping criterion can be

$$\max_i \left(\max_{\alpha_i^m < C_{\bar{y}_i}^m} -\nabla f(\alpha)_i^m - \min_{\alpha_i^m \leq C_{\bar{y}_i}^m} -\nabla f(\alpha)_i^m \right) < \epsilon \quad (38)$$

which is directly derived from (34) and (35). Note that since there are no lower bounds and all coefficients in (37a) are +1, the condition $\alpha_i^m > o_i^m, y_i^m = 1$ in (34) becomes $\alpha_i^m \leq C_{\bar{y}_i}^m$.

Now there is no lower bound on α_i^m , so our requirement that lower bounds must be larger than $-\infty$ seems to be violated. However, for this problem we can easily set a finite lower bound on α_i^m as follows. Using (37a)

$$\begin{aligned}\alpha_i^m &= -\sum_{m' \neq m} \alpha_i^{m'} \geq -\sum_{m' \neq m} C_{y_i}^{m'} \\ &= -C > -2C.\end{aligned}\quad (39)$$

Thus, (37) is still in the form of (33). Since the feasibility is always kept, α_i^m never reaches the lower bound $-2C$. Hence, the stopping criterion (35) still reduces to (38).

Implementations of decomposition methods for (37) have been discussed in [4] and [5]. Basically, the i th group of variables which has the maximal violation in (35) becomes the working set. Thus, the subproblem at the k th iteration is

$$\begin{aligned}\min_{\alpha_i} \quad & \frac{1}{2} K_{ii} \alpha_i^T \alpha_i + \left(\bar{e}_i + \sum_{j \neq i} K_{ij} (\alpha^k)_j \right)^T \alpha_i \\ \text{subject to} \quad & \sum_{m=1}^r \alpha_i^m = 0, \quad \alpha_i^m \leq C_{y_i}^m, \quad m = 1, \dots, r\end{aligned}\quad (40)$$

where $(\alpha^k)_j$ means $[(\alpha^k)_j^1, \dots, (\alpha^k)_j^r]^T$. Note that (40) is a very simple problem. In [4], two methods were proposed for it: one is an $O(r \log r)$ algorithm, while the other is an iterative procedure.

Since a whole group of variables $\alpha_1^1, \dots, \alpha_1^r$ is selected, Conditions 2'-4' of Theorem III.1 hold. If the sequence converges to an optimal solution, we have (36). In Section IV, we will prove the convergence of this decomposition method.

IV. CONVERGENCE OF A DECOMPOSITION METHOD FOR MULTICLASS SVM BY CRAMMER AND SINGER

The decomposition method mentioned in Section III for multiclass SVM is not in the category of decomposition methods considered in [8]. Hence, proofs in [8] cannot be directly used. However, since the working set selection as well as the subproblem are quite special where each time one of the l groups of variables is considered, we will show that the convergence proof is even simpler. First, we prove a simple lemma.

Lemma IV.1: Consider the following problem:

$$\begin{aligned}\min_s \quad & \frac{1}{2} s^T Q' s + p^T s \\ \text{subject to} \quad & y^T s = 0 \\ & o_i \leq s_i \leq u_i, \quad i = 1, \dots, l\end{aligned}\quad (41)$$

where $y_i = \pm 1, u_i \geq 0$, and $o_i \leq 0, i = 1, \dots, l$. Let $\min(\text{eig}(\cdot))$ be the smallest eigenvalue of a matrix. If there is $\sigma > 0$ such that $\min(\text{eig}(Q)) > \sigma$, then at an optimal solution s of (41)

$$\frac{1}{2} s^T Q' s + p^T s \leq -\frac{\sigma}{2} \|s\|^2.$$

Proof: From the KKT condition of (41), if s is an optimal solution, there is a b such that

$$\begin{aligned}(Qs)_i + p_i + by_i &= 0, & \text{if } o_i < s_i < u_i \\ (Qs)_i + p_i + by_i &\geq 0, & \text{if } s_i = o_i \\ (Qs)_i + p_i + by_i &\leq 0, & \text{if } s_i = u_i.\end{aligned}$$

Since $o_i \leq 0$ and $u_i \geq 0$

$$s^T Qs + p^T s + by^T s \leq 0.$$

With $y^T s = 0$

$$\frac{1}{2} s^T Qs + p^T s \leq -\frac{1}{2} s^T Qs \leq -\frac{\sigma}{2} \|s\|^2.$$

To use Lemma IV.1, we make the following assumption.

Assumption IV.2: There exists $\delta > 0$ such that for all $i = 1, \dots, l$, the kernel matrix K satisfies

$$K_{ii} \geq \delta.$$

From now on, we consider any one of convergent subsequences $\{\alpha^k\}, k \in \mathcal{K}$ and assume

$$\lim_{k \rightarrow \infty, k \in \mathcal{K}} \alpha^k = \bar{\alpha}.\quad (42)$$

We then have the following lemma.

Lemma IV.3: Assume $\{\alpha^k\}$ and $\bar{\alpha}$ are as defined in (42). Then for any given positive integer s , the sequence $\{\alpha^{k+s}\}, k \in \mathcal{K}$ converges to $\bar{\alpha}$.

Proof: If the i th group is selected as the working set, we have

$$f(\alpha) - f(\alpha^k) = \frac{1}{2} K_{ii} d^T d + p^T d$$

where

$$\begin{aligned}d &= [\alpha_i^1 - (\alpha^k)_i^1, \dots, \alpha_i^r - (\alpha^k)_i^r]^T \\ &= \alpha_i - (\alpha^k)_i \quad \text{and} \quad p = \bar{e}_i + \sum_{j=1}^l K_{ij} (\alpha^k)_j.\end{aligned}$$

Hence, an equivalent form of the subproblem (40) is

$$\begin{aligned}\min_d \quad & \frac{1}{2} K_{ii} d^T d + p^T d \\ \text{subject to} \quad & \sum_{m=1}^r d_m = 0 \\ & d^m \leq C_{y_i}^m - (\alpha^k)_i^m, \quad m = 1, \dots, r\end{aligned}\quad (43)$$

where d is the variable. Since α^k is a feasible point of (37), $C_{y_i}^m - (\alpha^k)_i^m \geq 0$.

As the smallest eigenvalue of the Hessian of (43) is K_{ii} and (43) is in the form of (41), from Assumption IV.2 and Lemma IV.1, we have

$$f(\alpha^{k+1}) - f(\alpha^k) \leq -\frac{\sigma}{2} \|\alpha^{k+1} - \alpha^k\|^2.\quad (44)$$

Next, we show that $\{f(\alpha^k)\}$ is a convergent sequence. First, we know that $\{f(\alpha^k)\}$ is decreasing. Using (39), the feasible region of (37) is a compact set, so $\lim_{k \rightarrow \infty} f(\alpha^k)$ exists and

$$\lim_{k \rightarrow \infty} f(\alpha^k) - f(\alpha^{k+1}) = 0. \quad (45)$$

Then, for the subsequence $\{\alpha^{k+1}\}, k \in \mathcal{K}$, from (44) and (45) we have

$$\begin{aligned} & \lim_{k \rightarrow \infty} \|\alpha^{k+1} - \bar{\alpha}\| \\ & \leq \lim_{k \rightarrow \infty} (\|\alpha^{k+1} - \alpha^k\| + \|\alpha^k - \bar{\alpha}\|) \\ & \leq \lim_{k \rightarrow \infty} \left(\sqrt{\frac{2}{\sigma}(f(\alpha^k) - f(\alpha^{k+1}))} + \|\alpha^k - \bar{\alpha}\| \right) = 0. \end{aligned}$$

Thus

$$\lim_{k \rightarrow \infty, k \in \mathcal{K}} \alpha^{k+1} = \bar{\alpha}.$$

From $\{\alpha^{k+1}\}$, we can prove $\lim_{k \rightarrow \infty, k \in \mathcal{K}} \alpha^{k+2} = \bar{\alpha}$ too. Therefore, $\lim_{k \rightarrow \infty, k \in \mathcal{K}} \alpha^{k+s} = \bar{\alpha}$ for any given s . ■

We then need a technical lemma.

Lemma IV.4: Assume $\{\alpha^k\}$ and $\bar{\alpha}$ are as defined in (42). If $m(\bar{\alpha})_i > M(\bar{\alpha})_i$, then after $k \in \mathcal{K}$ is large enough, $(\alpha^k)_i^m, m = 1, \dots, r$ are not changed.

Proof: Assume $(\alpha^k)_i, k \in \mathcal{K}$ is selected and changed infinitely many times. At any of these $(\alpha^k)_i$, we solve the subproblem (40) to obtain $(\alpha^{k+1})_i$. Now, consider the following problem with variable α_i

$$\begin{aligned} & \min_{\alpha_i} \frac{1}{2} K_{ii} \alpha_i^T \alpha_i + \left(\bar{e}_i + \sum_{j \neq i} K_{ij}(\bar{\alpha}_j) \right)^T \alpha_i \\ & \text{subject to} \quad \sum_{m=1}^r \alpha_i^m = 0 \\ & \quad \alpha_i^m \leq C_{\bar{y}_i}^m, \quad m = 1, \dots, r. \end{aligned} \quad (46)$$

Since $m(\bar{\alpha})_i > M(\bar{\alpha})_i$, $\bar{\alpha}_i$ is not an optimal solution of (46), we assume an optimal solution of (46) is $\tilde{\alpha}_i$.

Since α_i^{k+1} is an optimal solution of (40), we have

$$\begin{aligned} & \frac{1}{2} K_{ii} (\alpha^{k+1})_i^T (\alpha^{k+1})_i + \left(\bar{e}_i + \sum_{j \neq i} K_{ij}(\alpha^k)_j \right)^T (\alpha^{k+1})_i \\ & \leq \frac{1}{2} K_{ii} \tilde{\alpha}_i^T \tilde{\alpha}_i + \left(\bar{e}_i + \sum_{j \neq i} K_{ij}(\alpha^k)_j \right)^T \tilde{\alpha}_i. \end{aligned} \quad (47)$$

As $k \in \mathcal{K}$ goes to infinity, from Lemma IV.3, (47) becomes

$$\begin{aligned} & \frac{1}{2} K_{ii} \tilde{\alpha}_i^T \tilde{\alpha}_i + \left(\bar{e}_i + \sum_{j \neq i} K_{ij} \bar{\alpha}_j \right)^T \tilde{\alpha}_i \\ & \leq \frac{1}{2} K_{ii} \tilde{\alpha}_i^T \tilde{\alpha}_i + \left(\bar{e}_i + \sum_{j \neq i} K_{ij} \bar{\alpha}_j \right)^T \tilde{\alpha}_i \end{aligned}$$

which contradicts that $\bar{\alpha}_i$ is not an optimal solution of (46). Therefore, $(\alpha^k)_i$ is not selected after k is large enough. Hence, $(\alpha^k)_i$ remains the same. ■

The main result on the convergence is the following.

Theorem IV.5: Assume $\{\alpha^k\}$ is the sequence generated using the algorithm by Cramer and Singer. Under the Assumption IV.2, for any convergent subsequence of $\{\alpha^k\}$, its limit point is a global minimum of (37).

Proof: Using (39), we know that the feasible region of (37) is compact. Hence, $\{\alpha^k\}$ has convergent subsequences. Assume $\bar{\alpha}$ is the limit point of one convergent subsequence $\{\alpha^k\}, k \in \mathcal{K}$. If $\bar{\alpha}$ is not an optimal solution of (37), from the KKT condition (34), there are some groups j such that

$$m(\bar{\alpha})_j > M(\bar{\alpha})_j.$$

We define

$$\delta \equiv \min\{m(\bar{\alpha})_j - M(\bar{\alpha})_j \mid m(\bar{\alpha})_j - M(\bar{\alpha})_j > 0\}. \quad (48)$$

Using Lemma IV.4 and the continuity of $\nabla f(\alpha)$, we can consider all $k \in \mathcal{K}$ large enough such that the following three statements are valid.

- 1) Those $(\alpha^k)_j$ satisfying $m(\bar{\alpha})_j > M(\bar{\alpha})_j$ are not changed any more.
- 2) For all $i = 1, \dots, l, m = 1, \dots, r, k \leq k_1 \leq k + l$

$$|\nabla f(\alpha^{k_1})_i^m - \nabla f(\bar{\alpha})_i^m| < \delta/8. \quad (49)$$

- 3) For all $i = 1, \dots, l, m = 1, \dots, r, k \leq k_1, k_2 \leq k + l$

$$|\nabla f(\alpha^{k_1})_i^m - \nabla f(\alpha^{k_2})_i^m| < \delta/(8l). \quad (50)$$

In the rest of this proof, we will have a procedure to show that some group j with $m(\bar{\alpha})_j > M(\bar{\alpha})_j$ is still selected. This then contradicts Lemma IV.4.

Consider the k th iteration where a group i_1 is selected and modified. Thus

$$m(\alpha^k)_{i_1} > M(\alpha^k)_{i_1} \quad \text{but} \quad m(\alpha^{k+1})_{i_1} \leq M(\alpha^{k+1})_{i_1}. \quad (51)$$

Using (38), we assume that

$$-\nabla f(\alpha^{k+1})_{i_1}^{m_1} = m(\alpha^{k+1})_{i_1}. \quad (52)$$

Then, at the $(k+1)$ st iteration, if $i_2 \neq i_1$ is selected, then $(\alpha^{k+1})_{i_1} = (\alpha^{k+2})_{i_1}$ is not changed. Thus, if

$$-\nabla f(\alpha^{k+2})_{i_1}^{m_2} = m(\alpha^{k+2})_{i_1} \quad (53)$$

using (50), (52), (53), and the definition of $m(\alpha)$, we have

$$\begin{aligned} -\nabla f(\alpha^{k+1})_{i_1}^{m_1} - \delta/(8l) & \leq -\nabla f(\alpha^{k+2})_{i_1}^{m_1} \\ & \leq -\nabla f(\alpha^{k+2})_{i_1}^{m_2} \\ & \leq -\nabla f(\alpha^{k+1})_{i_1}^{m_2} + \delta/(8l) \\ & \leq -\nabla f(\alpha^{k+1})_{i_1}^{m_1} + \delta/(8l). \end{aligned}$$

Thus

$$\begin{aligned} & |-\nabla f(\alpha^{k+1})_{i_1}^{m_1} + \nabla f(\alpha^{k+2})_{i_1}^{m_2}| \\ & = |m(\alpha^{k+1})_{i_1} - m(\alpha^{k+2})_{i_1}| \leq \delta/(8l). \end{aligned} \quad (54)$$

Similarly

$$|M(\alpha^{k+1})_{i_1} - M(\alpha^{k+2})_{i_1}| \leq \delta/(8l)$$

so with (51)

$$m(\alpha^{k+2})_{i_1} \leq M(\alpha^{k+2})_{i_1} + \delta/(4l). \quad (55)$$

Thus, in the next l iterations, no matter i_1 is selected or not, we have

$$\begin{aligned} m(\alpha^{k+1})_{i_1} &\leq M(\alpha^{k+1})_{i_1} + \delta/4, \dots, m(\alpha^{k+l})_{i_1} \\ &\leq M(\alpha^{k+l})_{i_1} + \delta/4. \end{aligned} \quad (56)$$

On the other hand, from (48) and (49) and the fact that

$$\alpha_j^k = \dots = \alpha_j^{k+l} = \bar{\alpha}_j \quad (57)$$

some group j with $m(\bar{\alpha})_j > M(\bar{\alpha})_j$ satisfies

$$\begin{aligned} m(\alpha^{k+1})_j - M(\alpha^{k+1})_j &\geq m(\bar{\alpha})_j - M(\bar{\alpha})_j - \delta/4 \\ &\geq 3\delta/4, \dots, m(\alpha^{k+l})_j - M(\alpha^{k+l})_j \geq 3\delta/4. \end{aligned} \quad (58)$$

Note that the first inequality of (58) comes from a similar derivation of (55). For (55), in (54) the difference between $m(\alpha^{k+1})_{i_1}$ and $m(\alpha^{k+2})_{i_1}$ is estimated. For (58), since (57), we can directly measure the difference between $m(\alpha^{k+1})_j$ and $m(\bar{\alpha})_j$.

Hence, (56) and (58) imply that i_1 should not be selected in the $(k+1), \dots, (k+l)$ th iterations. Therefore, since there are l groups of variables, in l iterations, some group j with $m(\bar{\alpha})_j > M(\bar{\alpha})_j$ must be selected. This contradicts Lemma IV.4, which shows that this j th group of variables should not be selected.

Therefore, any limit point of $\{\alpha^k\}$ is a KKT point of (37). As (37) is a convex optimization problem, any limit point is a global minimum. ■

If the kernel matrix K is positive definite, $K \otimes I$ is also positive definite. Then, (37) is a strictly convex problem so there is a unique optimal solution. Thus, $\{\alpha^k\}$ is a globally convergent sequence if K is positive definite.

V. CONCLUSION AND DISCUSSION

Originally, we tried to prove Theorem II.1 by using as few conditions on the decomposition methods as possible. Surprisingly, we finally needed most properties of an existing working set selection. After finishing the convergence proof [8], an open question left is whether more flexible working set selections still lead to convergence. So far, unfortunately, in many scenarios, properties of a systematic working selection are always needed. This seems to hint that proving more generalized convergence may not be an easy task.

Next, we discuss the convergence proof for the formulation by Crammer and Singer. We can see that Assumption IV.2 is generally true. For example, for the polynomial kernel $K(x, y) = (x^T y)^d$, if all data are not zero vectors, $K_{ii} = (x_i^T x_i)^d > 0$. On the other hand, in [8] it assumes $\min_I(\min(\text{eig}(K_{II}))) > 0$, where I is any subset of $\{1, \dots, l\}$ with $|I| \leq q$, K_{II} is a square submatrix of K , and $\min(\text{eig}(\cdot))$ is the smallest eigenvalue of a matrix. We have to consider any I as there are no restrictions on the working set. On the other hand, the reason why Assumption IV.2 is simpler is that now in each iteration one of the l groups is selected where each group

has r variables $\alpha_i^1, \dots, \alpha_i^r$. Hence, the square submatrix of $K \otimes I$ is reduced to a small diagonal matrix $K_{ii} \otimes I$. Then all eigenvalues of $K_{ii} \otimes I$ are K_{ii} .

The proof in this paper also shows the importance of Lemma IV.3. For both proofs here and in [8], as we are not able to prove the global convergence of $\{\alpha^k\}$, instead we prove that $\{\alpha^{k+1}\}, k \in \mathcal{K}$ converges if $\{\alpha^k\}, k \in \mathcal{K}$ is a convergent subsequence. Then, this property is used to link several subproblems in subsequent iterations.

ACKNOWLEDGMENT

The author thanks four anonymous referees for their helpful comments.

REFERENCES

- [1] C.-C. Chang and C.-J. Lin. (2001) LIBSVM: A Library for Support Vector Machines. [Online]. Available: <http://www.csie.ntu.edu.tw/~cjlin/libsvm>
- [2] —, "Training ν -support vector classifiers: Theory and algorithms," *Neural Comput.*, vol. 13, no. 9, pp. 2119–2147, 2001.
- [3] C. Cortes and V. Vapnik, "Support-vector network," *Machine Learning*, vol. 20, pp. 273–297, 1995.
- [4] K. Crammer and Y. Singer, "On the algorithmic implementation of multiclass kernel-based vector machines," *J. Machine Learning R.*, vol. 2, pp. 265–292, 2001.
- [5] C.-W. Hsu and C.-J. Lin, "A comparison of methods for multiclass support vector machines," *IEEE Trans. Neural Networks*, vol. 13, pp. 415–425, Mar. 2002.
- [6] T. Joachims, *Advances in Kernel Methods—Support Vector Learning*, B. Schölkopf, C. J. C. Burges, and A. J. Smola, Eds. Cambridge, MA: MIT Press, 1998.
- [7] 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.
- [8] 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.
- [9] —, "Asymptotic convergence of an SMO algorithm without any assumptions," *IEEE Trans. Neural Networks*, vol. 13, pp. 248–250, Jan. 2002.
- [10] E. Osuna, R. Freund, and F. Girosi, "Training support vector machines: An application to face detection," in *Proc. CVPR'97*, 1997, pp. 130–136.
- [11] J. C. Platt, *Advances in Kernel Methods—Support Vector Learning*, B. Schölkopf, C. J. C. Burges, and A. J. Smola, Eds. Cambridge, MA: MIT Press, 1998.
- [12] B. Schölkopf, A. Smola, R. C. Williamson, and P. L. Bartlett, "New support vector algorithms," *Neural Comput.*, vol. 12, pp. 1207–1245, 2000.
- [13] V. Vapnik, *Statistical Learning Theory*. New York: Wiley, 1998.



Chih-Jen Lin (S'91–M'98) received the B.S. degree in mathematics from National Taiwan University, Taipei, Taiwan, in 1993. He received the M.S. and Ph.D. degrees from the Department of Industrial and Operations Engineering, University of Michigan, Ann Arbor, in 1997 and 1998, respectively.

Since September 1998, he has been an Assistant Professor in the Department of Computer Science and Information Engineering, National Taiwan University. His research interests include machine learning, numerical optimization, and various

applications of operations research.