

General Polynomial Time Decomposition Algorithms [★]

Nikolas List and Hans Ulrich Simon

Fakultät für Mathematik, Ruhr-Universität Bochum, 44780 Bochum, Germany
nlist@lmi.rub.de, simon@lmi.rub.de

Abstract. We present a general decomposition algorithm that is uniformly applicable to every (suitably normalized) instance of Convex Quadratic Optimization and efficiently approaches the optimal solution. The number of iterations required to be within ε of optimality grows linearly with $1/\varepsilon$ and quadratically with the number m of variables. The working set selection can be performed in polynomial time. If we restrict our considerations to instances of Convex Quadratic Optimization with at most k_0 equality constraints for some fixed constant k_0 plus some so-called box-constraints (conditions that hold for most variants of SVM-optimization), the working set is found in linear time. Our analysis builds on a generalization of the concept of rate certifying pairs that was introduced by Hush and Scovel. In order to extend their results to arbitrary instances of Convex Quadratic Optimization, we introduce the general notion of a rate certifying q -set. We improve on the results of Hush and Scovel [8] in several ways. First our result holds for Convex Quadratic Optimization whereas the results of Hush and Scovel are specialized to SVM-optimization. Second, we achieve a higher rate of convergence even for the special case of SVM-optimization (despite the generality of our approach). Third, our analysis is technically simpler.

1 Introduction

Support vector machines (SVMs) introduced by Vapnik and co-workers [1, 29] are a promising technique for classification, function approximation, and other key problems in statistical learning theory. In this paper, we consider the optimization problems that are induced by SVMs, which are special cases of Convex Quadratic Optimization.

The difficulty of solving problems of this kind is the density of the matrix that represents the “quadratic part” of the cost function. Thus, a prohibitive amount of memory is required to store the matrix and traditional optimization algorithms (such as Newton, for example) cannot be directly applied. Several authors have

[★] This work was supported in part by the IST Programme of the European Community, under the PASCAL Network of Excellence, IST-2002-506778. This publication only reflects the authors’ views. This work was furthermore supported by the Deutsche Forschungsgemeinschaft Grant SI 498/7-1.

proposed (different variants of) a decomposition method to overcome this difficulty [25, 9, 26, 27, 21, 22, 3, 11, 16, 4, 23, 12, 10, 17, 18, 13, 7, 14, 8, 20, 19, 5]. Given an instance of Convex Quadratic Optimization, this method keeps track of a current feasible solution which is iteratively improved. In each iteration the variable indices are split into a “working set” $I \subseteq \{1, \dots, m\}$ and its complement $J = \{1, \dots, m\} \setminus I$. Then, the simplified instance with the variables $x_i, i \in I$, is solved, thereby leaving the values for the remaining variables $x_j, j \in J$, unchanged. The success of the method depends in a quite sensitive manner on the policy for the selection of the working set I (whose size is typically bounded by a small constant). Ideally, the selection procedure should be computationally efficient and, at the same time, effective in the sense that the resulting sequence of feasible solutions converges (with high speed) to an optimal limit point.

Our results and their relation to previous work: Hush and Scovel considered a special SVM-optimization problem that we denote as SVO in our paper. They introduced the notion of an “ α -rate certifying pair” and showed that every decomposition algorithm for SVO that always inserts an α -rate certifying pair in its current working set comes within ε of optimality after $O(1/(\varepsilon\alpha^2))$ iterations. Building on a result of Chang, Hsu, and Lin [3], they presented furthermore an algorithm that constructs an $1/m^2$ -rate certifying pair in $O(m \log m)$ steps.¹ Combining these results, we see that the decomposition algorithm of Hush and Scovel for problem SVO is within ε of optimality after $O(m^4/\varepsilon)$ iterations.

In this paper we present an extension of (and an improvement on) this result. We first define the general notion of an α -rate certifying q -set and show (with a simplified analysis) that it basically fits the same purpose for Convex Quadratic Optimization (denoted as CQO in our paper) as the α -rate certifying pair for SVO, where the number of iterations needed to be within ε of optimality is proportional to $q/(\varepsilon\alpha^2)$. We present a general decomposition algorithm that is uniformly applicable to every (suitably normalized) instance of CQO. Given an instance with k equality constraints and m variables, it finds an $1/m$ -rate certifying $(k + 1)$ -set in polynomial time.² Combining these results, we are within ε of optimality after $O(km^2/\varepsilon)$ iterations of our decomposition algorithm. The SVM-optimization problem SVO (considered by Hush and Scovel) has only one equality constraint. Plugging in $k = 1$ in our general result, we arrive at an upper bound on the number of iterations that improves on the bound obtained by Hush and Scovel by factor m^2 . The analysis of Hush and Scovel in [8] builds on an earlier analysis of conditional gradient algorithms by Dunn [6]. For this part of the analysis, we will present simpler arguments.

There are some alternatives to the approach with rate certifying pairs. The most prominent one is the selection of the maximally KKT-violating pairs as

¹ Time bound $O(m \log m)$ can be improved to $O(m)$ by using the method from [28].

² Moreover, the algorithm can be implemented such as to find this set even in linear time when we restrict its application to instances of CQO with at most k_0 equality constraints for some fixed constant k_0 . If we restrict its application to SVO, we may use the the highly efficient method from [28].

implemented for example in SVM^{light} [9] or LIBSVM [4]. A paper by Lin [15] seems to imply the following quite strong result for SVO (although not stating it explicitly): decomposition algorithms following the approach of maximally KKT-violating pairs are within ε of optimality after only $O(\log 1/\varepsilon)$ iterations.³ However, the analysis is specialized to SVO. Furthermore it has to assume strict convexity of the objective function and some non-degeneracy conditions. The convergence rate is only given in terms of ε whereas the dependence on problem parameters (like, for example, m) is not clarified.

Another algorithm (related to but different from decomposition algorithms) is SimpleSVM [30] which tries to iteratively include the support vectors in the working set. Assuming strict convexity of the objective function, the authors of [30] claim a linear convergence of the method (but do neither give a complete proof nor exhibit the dependence on the various parameters explicitly). The main difference between SimpleSVM and decomposition algorithms is the size of the working set which can grow-up to the number of support vectors in the former case and is kept constant in the latter. Note that the number of support vectors is particularly large on noisy data.

There are many other papers about decomposition algorithms or related approaches that are noteworthy but cannot be mentioned properly in this abstract. The reader interested in finding more pointers to the relevant literature is referred to the full paper.

2 Preliminaries

We are mainly concerned with the problem “Convex Quadratic Optimization with box-constraints”. It is denoted simply as CQO in this paper and is formally given as follows:

Definition 1 (CQO). *An instance \mathcal{P} of CQO is given by*

$$\min_x f(x) = \frac{1}{2}x^\top Qx + w^\top x \text{ s.t. } Ax = b, l \leq x \leq r \text{ ,}$$

where

- $Q \in \mathbb{R}^{m \times m}$ is a symmetric positive semi-definite matrix over the reals and $w \in \mathbb{R}^m$, i.e., $f(x)$ is a convex quadratic cost function in m scalar variables,
- $A \in \mathbb{R}^{k \times m}$ and $b \in \mathbb{R}^k$ such that $Ax = b$ represents k linear equality constraints,
- $l, r \in \mathbb{R}^m$ and $l \leq x \leq r$ is the short-notation for the “box-constraints”

$$\forall i = 1, \dots, m : l_i \leq x_i \leq r_i \text{ .}$$

³ See [5] for a generalization of this result to similar but more general policies for working set selection.

In this paper, we will sometimes express $f(x')$ by means of the Taylor-expansion around x :

$$f(x') = f(x) + \nabla f(x)^\top (x' - x) + \frac{1}{2} (x' - x)^\top Q (x' - x) ,$$

where $\nabla f(x) = Qx + w$. For $d := x - x'$, this can be rewritten as follows:

$$f(x) - f(x') = \nabla f(x)^\top d - \frac{1}{2} d^\top Q d . \quad (1)$$

Note that $d^\top Q d \geq 0$ because Q is positive semi-definite.

In the sequel,

$$R(\mathcal{P}) = \{x \in \mathbb{R}^m \mid Ax = b, l \leq x \leq r\}$$

denotes the compact and convex set of feasible points for \mathcal{P} . The well-known first-order condition for convex function optimization⁴ (valid for an arbitrary convex cost function) states that x is an optimal feasible solution iff

$$\forall x' \in R(\mathcal{P}) : \nabla f(x)^\top (x' - x) \geq 0 .$$

We briefly note that any instance of the *general* convex quadratic optimization problem with cost function $f(x)$, linear equality constraints, linear inequality constraints (not necessarily in the form of box-constraints) and a compact region of feasible points can be transformed into an equivalent instance of CQO because we may convert the linear inequalities into linear equations by introducing non-negative slack variables. By the compactness of the region of feasible points, we may also put a suitable upper bound on each slack variable such that finally all linear inequalities take the form of box-constraints.

We now define a subproblem of CQO, denoted as SVO in this paper, that is actually one of the most well studied SVM-optimization problems:

Definition 2 (SVO). *An instance \mathcal{P}_0 of SVO is given by*

$$\min_x f(x) \text{ s.t. } y^\top x = 0 , l \leq x \leq r ,$$

where $f(x), l, r$ are understood as in Definition 1 and $y \in \{-1, 1\}^m$ is a vector whose components represent binary classification labels.

The main difference between SVO and the general problem CQO is that SVO has only a single equality constraint. Furthermore, this equality constraint is of a special form.

We are now prepared to introduce (informally) the notion of “decomposition algorithms”. A *decomposition algorithm for CQO* with working sets of size at most q (where we allow that q depends on k) proceeds iteratively as follows: given an instance \mathcal{P} of CQO and a feasible solution $x \in R(\mathcal{P})$ (chosen arbitrarily in the beginning), a so-called working set $I \subseteq \{1, \dots, m\}$ of size at most q is selected. Then x is updated by the optimal solution for the simplified instance with variables $x_i, i \in I$ (leaving the values x_j with $j \notin I$ unchanged). *Decomposition algorithms for SVO* are defined analogously. The policy for working set selection is a critical issue that we discuss in the next sections.

⁴ see for example [2].

Notational Conventions:

- The parameters m, k, f, A, y, b, l, r are consistently used in this paper as the components of an instance of CQO or SVO. Similarly, parameter q (possibly dependent on k) always represents the (maximal) size of the working set.
- L_{max} and S_{max} are two more parameters that we associate with an instance of CQO or SVO (where L_{max} depends also on q). L_{max} denotes the largest among the eigenvalues of all the principal $(q \times q)$ -submatrices of Q . S_{max} denotes the maximum side length of the box spanned by l and r , i.e., $S_{max} := \max_{1 \leq i \leq m} (r_i - l_i)$.
- For a decomposition algorithm \mathcal{A} , we denote the current feasible solution obtained after n iterations by x^n (such that x^0 is the feasible solution \mathcal{A} starts with). The optimal feasible solution is denoted by x^* . Then

$$\Delta_n := f(x^n) - f(x^*) \quad (2)$$

denotes the difference between the value of the current solution and the optimal value.

3 Rate Certifying Sets and the Main Theorem

In section 3.1, we recall the concept of rate certifying pairs. In section 3.2, we present the new notion of rate certifying sets and state our main result, whose proof is given in sections 3.3 and 3.4.

3.1 Rate Certifying Pairs

We consider again the problem SVO from Definition 2 along with a problem instance \mathcal{P}_0 . Let $x \in R(\mathcal{P}_0)$ be a feasible solution and $x^* \in R(\mathcal{P}_0)$ an optimal feasible solution. In the sequel, we will often use the following first-order approximation of the maximal distance (with regard to the value of the objective function) between a given point x and any other feasible solution x' :

$$\sigma(x) := \sup_{x' \in R(\mathcal{P}_0)} \nabla f(x)^\top (x - x') .$$

As already noted by Hush and Scovel [8], the following holds:⁵

$$f(x) - f(x^*) \leq \nabla f(x)^\top (x - x^*) \leq \sigma(x) . \quad (3)$$

In other words, $f(x)$ is always within $\sigma(x)$ of optimality. Note that $\sigma(x^*) = 0$ (which immediately follows from the first-order optimality condition). Thus, $\sigma(x) > 0$ if and only if x is suboptimal.

Since we are dealing with working sets whose size is bounded by a (small) parameter q , it is natural to restrict the range of x' to feasible solutions that

⁵ The first inequality follows from (1) and the positive semi-definiteness of Q ; the second one is trivial.

differ from x in at most q coordinates. For $q = 2$, this leads to the following definition:

$$\sigma(x|i_1, i_2) := \sup_{x' \in R(\mathcal{P}_0): x'_i = x_i \text{ for } i \neq i_1, i_2} \nabla f(x)^\top (x - x') .$$

The following notion is crucial: (i_1, i_2) is called an α -rate certifying pair for x if

$$\sigma(x|i_1, i_2) \geq \alpha(f(x) - f(x^*)) .$$

Let α be a function in m with strictly positive values. An α -rate certifying algorithm is a decomposition algorithm for SVO that, for every m and every input instance \mathcal{P}_0 with m variables, always includes an $\alpha(m)$ -rate certifying pair in the current working set. As mentioned already in the introduction, the main results in [8] are as follows:

Theorem 1 (Hush and Scovel [8]).

1. Let \mathcal{A} be an α -rate certifying algorithm. Consider any instance \mathcal{P}_0 of SVO with, say, m variables. Let L_{max} and S_{max} be the quantities associated with \mathcal{P}_0 .⁶ For sake of brevity, let $\alpha = \alpha(m)$. Then, \mathcal{A} is within ε of optimality after

$$1 + \frac{2 \max\{1, 2S_{max}^2\}}{\alpha} \left(\frac{\max\{1, \alpha \Delta_0 / L_{max}\} L_{max}}{\alpha \varepsilon} - 1 \right) = \\ O \left(\frac{L_{max}(1 + S_{max})^2}{\alpha^2 \varepsilon} + \frac{\Delta_0(1 + S_{max})^2}{\alpha \varepsilon} \right)$$

iterations.

2. For function α given by $\alpha(m) = 1/m^2$, there exists an α -rate certifying algorithm. It constructs a working set (given \mathcal{P}_0 and a suboptimal feasible solution x) in $O(m \log m)$ steps (or in $O(m)$ steps when the method from [28] is applied). Furthermore, it is within ε of optimality after

$$O \left(\frac{L_{max}(1 + S_{max})^2 m^4}{\varepsilon} + \frac{\Delta_0(1 + S_{max})^2 m^2}{\varepsilon} \right)$$

iterations.

3.2 Rate Certifying Sets

The definition of $\sigma(x)$ is easily extended to any instance \mathcal{P} of CQO:

$$\sigma(x) := \sup_{x' \in R(\mathcal{P})} \nabla f(x)^\top (x - x') . \quad (4)$$

Clearly, inequality (3) and the subsequent comments are still valid without any change. However, since CQO deals with several equality constraints, one can in

⁶ See our notational conventions in section 2 for a definition (where $q = 2$).

general not expect to find rate certifying *pairs*. Instead, the following general definition for $I \subseteq \{1, \dots, m\}$ will prove useful:

$$\sigma(x|I) := \sup_{x' \in R(\mathcal{P}): x'_i = x_i \text{ for } i \notin I} \nabla f(x)^\top (x - x') .$$

I is called an α -rate certifying q -set if $|I| \leq q$ and

$$\sigma(x|I) \geq \alpha(f(x) - f(x^*)) . \quad (5)$$

Let α be a function in m with strictly positive values and let \mathbf{q} be a function in k whose values are strictly positive integers. An (α, \mathbf{q}) -rate certifying algorithm is a decomposition algorithm for CQO that, for every m, k and any problem instance \mathcal{P} with m variables and k equality constraints, includes an $\alpha(m)$ -rate certifying $\mathbf{q}(k)$ -set in the current working set.⁷ With these notations, the following holds:

Theorem 2. 1. Let \mathcal{A} be an (α, \mathbf{q}) -rate certifying algorithm. Consider an instance \mathcal{P} of CQO with, say, m variables and k equality constraints. For sake of brevity, let $\alpha = \alpha(m)$, $q = \mathbf{q}(k)$, and let L_{max}, S_{max} be the quantities associated with \mathcal{P} and q . Then, \mathcal{A} is within ε of optimality after

$$\left\lceil \frac{2qL_{max}S_{max}^2}{\alpha^2\varepsilon} \right\rceil + \max \left\{ 0, \left\lceil \frac{2}{\alpha} \ln \left(\frac{\Delta_0}{\varepsilon} \right) \right\rceil \right\} \quad (6)$$

iterations. Moreover, if $qL_{max}S_{max}^2 \leq \varepsilon\alpha$, then $\max\{0, \lceil 2 \ln(\Delta_0/\varepsilon)/\alpha \rceil\}$ iterations (the second term in (6)) are enough.

2. For functions α, \mathbf{q} given by $\alpha(m) = 1/m$ and $\mathbf{q}(k) = k + 1$, there exists an (α, \mathbf{q}) -rate certifying algorithm \mathcal{A} . It constructs a working set (given \mathcal{P} and a suboptimal feasible solution x) in polynomial time. Moreover, if we restrict its application to instances of CQO with at most k_0 equality constraints for some fixed constant k_0 , there is a linear time bound for the construction of the working set.⁸

3. The algorithm \mathcal{A} from the preceding statement is within ε of optimality after

$$\left\lceil \frac{2(k+1)m^2L_{max}S_{max}^2}{\varepsilon} \right\rceil + \max \left\{ 0, \left\lceil 2m \ln \left(\frac{\Delta_0}{\varepsilon} \right) \right\rceil \right\}$$

iterations. Moreover, if $(k+1)L_{max}S_{max}^2 \leq \varepsilon/m$, then $\max\{0, \lceil 2m \ln(\Delta_0/\varepsilon) \rceil\}$ iterations are enough.

Clearly, the third statement in Theorem 2 follows directly from the first two statements (which will be proven in subsections 3.3 and 3.4, respectively).

A few comments on Theorem 2 are in place here. One might be tempted to think that an (α, \mathbf{q}) -rate certifying algorithm decreases (an upper bound on)

⁷ Finding an initial feasible solution in the beginning is equivalent to solving a standard LP. For an instance of an SVM-optimization problem an initial guess usually is the zero vector.

⁸ If we restrict its application to instances of SVO, we may use the the highly efficient method from [28].

the distance between the current feasible solution and the best feasible solution (with regard to the objective value) roughly by factor $1 - \alpha$ (for $\alpha := \alpha(m)$). If such a “contraction” took place, we would be within ε of optimality after only $O(\log(1/\varepsilon)/\alpha)$ iterations. This is however spurious thinking because the σ -function is *not* concerned with this distance itself but rather with a first-order approximation of it. We will see in the proof of Theorem 2 that a run of an (α, \mathbf{q}) -rate certifying algorithm can be decomposed into two phases. As long as the distance from the optimal value is large in comparison to (an upper bound on) the second order terms (phase 1), a contraction by factor $1 - \alpha/2$ takes place. However when we come closer to the optimal value (phase 2), the effect of the neglected second order terms becomes more significant and the convergence slows down (at least within our perhaps somewhat pessimistic analysis). Phase 1 leads to the term $\max\{0, \lceil 2 \ln(\Delta_0/\varepsilon)/\alpha \rceil\}$ in (6) whereas phase 2 leads to the term $\lceil (2qL_{max}S_{max}^2)/(\alpha^2\varepsilon) \rceil$.

3.3 Proof of the 1st Statement in Theorem 2

We will use the notation introduced in Theorem 2 and consider an arbitrary iteration of the (α, \mathbf{q}) -rate certifying algorithm \mathcal{A} when it is applied on input \mathcal{P} . To this end, let x denote a feasible but suboptimal solution, and let x^* be the optimal feasible solution. Let I be the subset of the working set that satisfies $|I| \leq q$ and (5). Let x' be a feasible solution that satisfies $x_i = x'_i$ for every $i \notin I$ and

$$\sigma(x|I) = \nabla f(x)^\top (x - x') = \nabla f(x)^\top d , \quad (7)$$

where $d = x - x'$. Combining (5) with (7), we get

$$\nabla f(x)^\top d \geq \alpha(f(x) - f(x^*)) = \alpha\Delta , \quad (8)$$

where $\Delta = f(x) - f(x^*) \geq 0$. For some parameter $0 \leq \lambda \leq 1$ (suitably chosen later), consider the feasible solution

$$x'' = x - \lambda d$$

on the line segment between x and x' . Taylor-expansion around x allows to relate $f(x)$ and $f(x'')$ as follows:

$$f(x) - f(x'') = \lambda \nabla f(x)^\top d - \lambda^2 \frac{1}{2} d^\top Q d .$$

Note that x'' (like x') satisfies $x''_i = x_i$ for every $i \notin I$. Thus, $x'' = x - \lambda d$ is a feasible solution that coincides with x outside the current working set. Thus, \mathcal{A} (finding the optimal feasible solution that coincides with x outside the working set) achieves in the next iteration a “cost reduction” of at least $f(x) - f(x'')$.⁹

⁹ “Achieving a cost reduction of a in the next iteration” means that the next iteration decreases the distance between the value of the current feasible solution and the value of the optimal feasible solution by at least a .

x'' depends on the parameter $0 \leq \lambda \leq 1$. In the sequel, we tune parameter λ such as to obtain a “large” cost reduction. To be on the safe side, we will however perform a “pessimistic analysis” where we substitute worst case bounds for $\nabla f(x)^\top d$ and $d^\top Qd$ respectively. Clearly

$$\max_{0 \leq \lambda \leq 1} \left(\lambda \nabla f(x)^\top d - \lambda^2 \frac{1}{2} d^\top Qd \right) \geq \max_{0 \leq \lambda \leq 1} \left(\lambda B - \frac{1}{2} \lambda^2 B' \right)$$

for any lower bound B on $\nabla f(x)^\top d$ and any upper bound B' on $d^\top Qd$. According to (8), $\alpha\Delta$ can serve as a lower bound on $\nabla f(x)^\top d$. It is easily seen that the following parameter U can serve as an upper bound on $d^\top Qd$:

$$U := qL_{max}S_{max}^2 \geq d^\top Qd . \quad (9)$$

This immediately follows from the definition of L_{max} and S_{max} and from the fact that d has at most q non-zero components. We conclude from this discussion that, for every $0 \leq \lambda \leq 1$, \mathcal{A} achieves in the next iteration a cost reduction of at least

$$h(\lambda) := \lambda\alpha\Delta - \frac{1}{2}\lambda^2U .$$

It is easily seen that function $h(\lambda)$ is maximized by setting

$$\lambda := \begin{cases} 1 & \text{if } \Delta > U/\alpha \\ \alpha\Delta/U & \text{if } \Delta \leq U/\alpha \end{cases} .$$

Case 1 $\Delta > U/\alpha$. Then \mathcal{A} achieves a cost reduction of at least

$$h(1) = \alpha\Delta - \frac{1}{2}U > \alpha\Delta/2 .$$

Thus the difference $\Delta = f(x) - f(x^*)$ will shrink after the next iteration to

$$\Delta - \frac{\alpha}{2}\Delta = \Delta \left(1 - \frac{\alpha}{2} \right) ,$$

or to a smaller value, which is a proper contraction.

Case 2 $\Delta \leq U/\alpha$. Then \mathcal{A} achieves a cost reduction of at least

$$h\left(\frac{\alpha\Delta}{U}\right) = \frac{\alpha^2\Delta^2}{2U} = \gamma\Delta^2 ,$$

where

$$\gamma := \frac{\alpha^2}{2U} . \quad (10)$$

Thus, the difference $\Delta = f(x) - f(x^*)$ will shrink after the next iteration to

$$\Delta - \gamma\Delta^2$$

or to a smaller value.

Recall from (2) that sequence $(\Delta_n)_{n \geq 0}$ keeps track of the difference between the value of the current feasible solution and the value of the optimal solution. In view of the two cases described above, our pessimistic analysis obviously runs through two phases:

Phase 1 As long as $\Delta_n > U/\alpha$, we calculate with cost reduction $\alpha\Delta_n/2$.

Phase 2 As soon as $\Delta_n \leq U/\alpha$, we calculate with cost reduction $\gamma\Delta_n^2$.

Let us first assume that $\varepsilon < U/\alpha$ (and postpone the case $\varepsilon \geq U/\alpha$ to the end of this subsection). The number of iterations in phase 1 is not larger than the smallest $n_0 \geq 0$ that satisfies the second inequality in

$$\Delta_0 \left(1 - \frac{\alpha}{2}\right)^{n_0} < \Delta_0 e^{-n_0\alpha/2} \leq \varepsilon < \frac{U}{\alpha} ,$$

i.e.,

$$n_0 := \max \left\{ 0, \left\lceil \frac{2}{\alpha} \ln \left(\frac{\Delta_0}{\varepsilon} \right) \right\rceil \right\} . \quad (11)$$

In phase 2, $(\Delta_n)_{n \geq n_0}$ evolves according to

$$\Delta_{n+1} \leq \Delta_n - \gamma\Delta_n^2 = \Delta_n(1 - \gamma\Delta_n) . \quad (12)$$

Recall that $\Delta_i = f(x^i) - f(x^*) \geq 0$ for every i . As for the iterations considered in phase 2 within our analysis, we can make the stronger (pessimistic) assumption $\Delta_i > 0$.¹⁰ Following Dunn [6], we can therefore consider the reciprocals $\delta_n := 1/\Delta_n$. Note that (12) and $\Delta_{n+1} > 0$ imply that $0 \leq \gamma\Delta_n < 1$ and so for each $n \geq n_0$ the following relation holds:

$$\delta_{n+1} - \delta_n \geq \frac{1}{\Delta_n(1 - \gamma\Delta_n)} - \frac{1}{\Delta_n} = \frac{\gamma}{1 - \gamma\Delta_n} \geq \gamma .$$

Therefore

$$\delta_n = \delta_{n_0} + \sum_{j=n_0}^{n-1} (\delta_{j+1} - \delta_j) \geq \gamma(n - n_0)$$

and consequently

$$\Delta_n = \frac{1}{\delta_n} \leq \frac{1}{\gamma(n - n_0)} .$$

It follows that $\Delta_n \leq \varepsilon$ after

$$n - n_0 := \left\lceil \frac{1}{\gamma\varepsilon} \right\rceil$$

iterations in phase 2. Thus the total number of iterations in both phases needed to be within ε of optimality is bounded by

$$n_0 + \left\lceil \frac{1}{\gamma\varepsilon} \right\rceil \stackrel{(10),(11)}{\leq} \left\lceil \frac{2U}{\alpha^2\varepsilon} \right\rceil + \max \left\{ 0, \left\lceil \frac{2}{\alpha} \ln \left(\frac{\Delta_0}{\varepsilon} \right) \right\rceil \right\}$$

¹⁰ Otherwise phase 2 ends with the optimal solution even earlier.

iterations. Plugging in the definition of U from (9), we obtain (6).

Let us now finally discuss the case $U = qL_{max}S_{max}^2 \leq \varepsilon\alpha$. Since $\varepsilon \geq U/\alpha$, we come within ε of optimality during phase 1. The number of iterations required for this is the smallest n_0 that satisfies

$$\Delta_0 \left(1 - \frac{\alpha}{2}\right)^{n_0} < \Delta_0 e^{-n_0\alpha/2} \leq \varepsilon .$$

Thus, we are within ε of optimality after

$$\max \left\{ 0, \left\lceil \frac{2}{\alpha} \ln \left(\frac{\Delta_0}{\varepsilon} \right) \right\rceil \right\}$$

iterations. This completes the proof of the first statement in Theorem 2.

3.4 Proof of the 2nd Statement in Theorem 2

We first give a short outline of the proof. Let x be a feasible but suboptimal solution for \mathcal{P} . According to (3), it is sufficient to efficiently construct a working set I such that $|I| \leq k + 1$ and

$$\sigma(x|I) \geq \frac{1}{m}\sigma(x) . \quad (13)$$

To this end, we will proceed as follows. We consider auxiliary instances $\mathcal{P}_x, \mathcal{P}'_x, \mathcal{P}''_x$ of the Linear Programming Problem (denoted as LP in the sequel). The optimal values for \mathcal{P}_x and \mathcal{P}'_x , are both shown to coincide with $\sigma(x)$. From \mathcal{P}'_x , we derive (basically by aggregating several equality constraints into a single-one) the instance \mathcal{P}''_x , whose optimal basic solution will represent a working set I of size at most $k + 1$. A comparison of the three problem instances will finally reveal that I satisfies (13).

We now move on to the technical implementation of this plan. Recall from (4) that

$$\sigma(x) = \sup_{x' \in R(\mathcal{P})} \nabla f(x)^\top (x - x') = \nabla f(x)^\top d ,$$

where $d = x - x'$. Thus $\sigma(x)$ is the optimal value of the following instance \mathcal{P}_x of LP:

$$\max_d \nabla f(x)^\top d \text{ s.t. } Ad = \mathbf{0}, l \leq x - d \leq r .$$

We set

$$\mu^+ := x - l \text{ and } \mu^- := r - x$$

and split d into positive and non-positive components d^+ and d^- respectively:

$$d = d^+ - d^-, \quad d^+, d^- \geq \mathbf{0}, \quad \forall i = 1, \dots, m : d_i^+ d_i^- = 0 . \quad (14)$$

With these notations, $\sigma(x)$ also coincides with the optimal value of the following instance \mathcal{P}'_x of LP:

$$\max_{d^+, d^-} \begin{pmatrix} \nabla f(x) \\ -\nabla f(x) \end{pmatrix}^\top \begin{pmatrix} d^+ \\ d^- \end{pmatrix}$$

subject to

$$\begin{aligned} [A, -A] \begin{pmatrix} d^+ \\ d^- \end{pmatrix} &= \mathbf{0} \\ \mathbf{0} &\leq d^+ \leq \mu^+ , \quad \mathbf{0} \leq d^- \leq \mu^- \end{aligned}$$

The third instance \mathcal{P}_x'' of LP that we consider has an additional slack variable ξ and is given as follows:

$$\max_{d^+, d^-, \xi} \begin{pmatrix} \nabla f(x) \\ -\nabla f(x) \end{pmatrix}^\top \begin{pmatrix} d^+ \\ d^- \end{pmatrix}$$

subject to

$$\forall i = 1, \dots, m : \mu_i^- = 0 \Rightarrow d_i^- = 0 , \quad \mu_i^+ = 0 \Rightarrow d_i^+ = 0 \quad (15)$$

$$[A, -A] \begin{pmatrix} d^+ \\ d^- \end{pmatrix} = \mathbf{0} \quad (16)$$

$$\sum_{i: \mu_i^+ > 0} \frac{1}{\mu_i^+} d_i^+ + \sum_{i: \mu_i^- > 0} \frac{1}{\mu_i^-} d_i^- + \xi = 1 \quad (17)$$

$$d^+, d^- \geq \mathbf{0} , \quad \xi \geq 0 \quad (18)$$

We briefly note that we do not have to count the equality constraints in (15) because we may simply remove the variables that are set to zero from the problem instance. Recall that matrix A represents k equality constraints. Thus, \mathcal{P}_x'' is a linear program in canonical form with $k + 1$ equality constraints. Its optimal basic feasible solution has therefore at most $k + 1$ non-zero components. The following observations (where d, d^+, d^- are related according to (14)) are easy to verify:

1. If $\begin{pmatrix} d^+ \\ d^- \end{pmatrix}$ is a feasible solution for \mathcal{P}'_x with value p , then $\frac{1}{m} \begin{pmatrix} d^+ \\ d^- \end{pmatrix}$ is a feasible solution for \mathcal{P}''_x with value p/m .
2. If $\begin{pmatrix} d^+ \\ d^- \end{pmatrix}$ is a feasible solution for \mathcal{P}''_x with value p , then $\begin{pmatrix} d^+ \\ d^- \end{pmatrix}$ is also a feasible solution for \mathcal{P}'_x with value p .

Recall that $\sigma(x)$ is the value of the optimal solution for \mathcal{P}'_x . We may conclude from our observations that the value of the optimal solution for \mathcal{P}''_x , say $\sigma'(x)$, satisfies $\sigma'(x) \geq \sigma(x)/m$. Now consider an optimal basic feasible solution $\begin{pmatrix} d^+ \\ d^- \end{pmatrix}$ for \mathcal{P}''_x with value $\sigma'(x)$. Let

$$I := \{i \in \{1, \dots, m\} \mid d_i^+ \neq 0 \text{ or } d_i^- \neq 0\} .$$

Clearly, $|I| \leq k + 1$. Since $d = d^+ - d^-$ is a feasible solution for \mathcal{P} (still with value $\sigma'(x)$) that differs from x only in coordinates from I , we may conclude that

$$\sigma(x|I) \geq \sigma'(x) \geq \frac{1}{m} \sigma(x) .$$

In other words, working set I satisfies (13). The time required to compute I is dominated by the time required to solve the linear program \mathcal{P}_x'' . This can be done in polynomial time. Since a linear program with a constant number of variables (or a linear program in standard form with a constant number of equality constraints¹¹) can be solved in linear time [24], the proof for the 2nd statement of Theorem 2 is completed.

4 Conclusions and Open Problems

We have presented an analysis of a decomposition algorithm that leads to the up-to-date strongest theoretical performance guarantees within the “rate certifying pair” approach. Our analysis holds uniformly for any instance of Convex Quadratic Optimization (with box-constraints) and certainly covers most of the variants of SVM-optimization. As explained in the introduction already, there are competing approaches like, for example, the approach based on maximally KKT-violating pairs or approaches based on an iterative inclusion of support vectors. As should become clear from the introduction, none of these approaches beats the other-ones in all respects (uniform analysis for a broad variety of problems, high speed of convergence, efficient working set selection). It remains an object of future research to gain more insight (theoretically and empirically) into the (perhaps complementary) strength and weakness of the various approaches such that their combined power can be exploited to full extent.

References

1. Bernhard E. Boser, Isabelle M. Guyon, and Vladimir N. Vapnik. A training algorithm for optimal margin classifiers. In *Proceedings of the 5th Annual ACM Workshop on Computational Learning Theory*, pages 144–152, 1992.
2. Steven Boyd and Lieven Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
3. Chih-Chung Chang, Chih-Wei Hsu, and Chih-Jen Lin. The analysis of decomposition methods for support vector machines. *IEEE Transactions on Neural Networks*, 11(4):248–250, 2000.
4. Chih-Chung Chang and Chih-Jen Lin. *LIBSVM: A library for support vector machines*, 2001. Available from <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
5. Pai-Hsuen Chen, Rong-En Fan, and Chih-Jen Lin. A study on SMO-type decomposition methods for support vector machines, 2005. Available from <http://www.csie.ntu.edu.tw/~cjlin/papers/generalSMO.pdf>.
6. J. Dunn. Rates of convergence for conditional gradient algorithms near singular and non-singular extremals. *SIAM J. Control and Optimization*, 17(2):187–211, 1979.
7. Chih-Wei Hsu and Chih-Jen Lin. A simple decomposition method for support vector machines. *Machine Learning*, 46(1–3):291–314, 2002.
8. Don Hush and Clint Scovel. Polynomial-time decomposition algorithms for support vector machines. *Machine Learning*, 51(1):51–71, 2003.

¹¹ such that the dual linear program has constant number of variables

9. Thorsten Joachims. Making large scale SVM learning practical. In Bernhard Schölkopf, Christopher J. C. Burges, and Alexander J. Smola, editors, *Advances in Kernel Methods—Support Vector Learning*, pages 169–184. MIT Press, 1998.
10. S. Sathiya Keerthi and E. G. Gilbert. Convergence of a generalized SMO algorithm for SVM classifier design. *Machine Learning*, 46(1–3):351–360, 2002.
11. S. Sathiya Keerthi, Shirish Krishnaji Shevade, Chiranjib Bhattacharyya, and K. R. K. Murthy. Improvements to SMO algorithm for SVM regression. *IEEE Transactions on Neural Networks*, 11(5):1188–1193, 2000.
12. S. Sathiya Keerthi, Shirish Krishnaji Shevade, Chiranjib Bhattacharyya, and K. R. K. Murthy. Improvements to Platt’s SMO algorithm for SVM classifier design. *Neural Computation*, 13(3):637–649, 2001.
13. Pavel Laskov. Feasible direction decomposition algorithms for training support vector machines. *Machine Learning*, 46(1–3):315–349, 2002.
14. Shuo-Peng Liao, Hsuan-Tien Lin, and Chih-Jen Lin. A note on the decomposition methods for support vector regression. *Neural Computation*, 14(6):1267–1281, 2002.
15. Chih-Jen Lin. Linear convergence of a decomposition method for support vector machines, 2001. Available from <http://www.csie.ntu.edu.tw/~cjlin/papers/linearconv.pdf>.
16. Chih-Jen Lin. On the convergence of the decomposition method for support vector machines. *IEEE Transactions on Neural Networks*, 12(6):1288–1298, 2001.
17. Chih-Jen Lin. Asymptotic convergence of an SMO algorithm without any assumptions. *IEEE Transactions on Neural Networks*, 13(1):248–250, 2002.
18. Chih-Jen Lin. A formal analysis of stopping criteria of decomposition methods for support vector machines. *IEEE Transactions on Neural Networks*, 13(5):1045–1052, 2002.
19. Nikolas List. Convergence of a generalized gradient selection approach for the decomposition method. In *Proceedings of the 15th International Conference on Algorithmic Learning Theory*, pages 338–349, 2004.
20. Nikolas List and Hans Ulrich Simon. A general convergence theorem for the decomposition method. In *Proceedings of the 17th Annual Conference on Computational Learning Theory*, pages 363–377, 2004.
21. Olvi L. Mangasarian and David R. Musicant. Successive overrelaxation for support vector machines. *IEEE Transactions on Neural Networks*, 10(5):1032–1037, 1999.
22. Olvi L. Mangasarian and David R. Musicant. Active support vector machine classification. In *Advances in Neural Information Processing Systems 12*, pages 577–583. MIT Press, 2000.
23. Olvi L. Mangasarian and David R. Musicant. Lagrangian support vector machines. *Journal of Machine Learning Research*, 1:161–177, 2001.
24. Nimrod Megiddo. Linear programming in linear time when the dimension is fixed. *Journal of the Association on Computing Machinery*, 31(1):114–127, 1984.
25. Edgar E. Osuna, Robert Freund, and Federico Girosi. Training support vector machines: an application to face detection. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, pages 130–136, 1997.
26. John C. Platt. Fast training of support vector machines using sequential minimal optimization. In Bernhard Schölkopf, Christopher J. C. Burges, and Alexander J. Smola, editors, *Advances in Kernel Methods—Support Vector Learning*, pages 185–208. MIT Press, 1998.
27. Craig Saunders, Mark O. Stitson, Jason Weston, Leon Bottou, Bernhard Schölkopf, and Alexander J. Smola. Support vector machine reference manual. Technical Report CSD-TR-98-03, Royal Holloway, University of London, Egham, UK, 1998.

28. Hans Ulrich Simon. On the complexity of working set selection. In *Proceedings of the 15th International Conference on Algorithmic Learning Theory*, pages 324–337, 2004.
29. Vladimir Vapnik. *Statistical Learning Theory*. Wiley Series on Adaptive and Learning Systems for Signal Processing, Communications, and Control. John Wiley & Sons, 1998.
30. S. V. N. Vishwanthan, Alexander J. Smola, and M. Narasimha Murty. SimpleSVM. In *Proceedings of the 20th International Conference on Machine Learning*, 2003.