



Parallel Variable Distribution for Constrained Optimization*

C.A. SAGASTIZÁBAL

sagastiz@impa.br

*INRIA-Rocquencourt, BP 105, 78153 Le Chesnay, France; Instituto de Matemática Pura e Aplicada,
Estrada Dona Castorina 110, Jardim Botânico, Rio de Janeiro, RJ 22460-320, Brazil*

M.V. SOLODOV

solodov@impa.br

*Instituto de Matemática Pura e Aplicada, Estrada Dona Castorina 110, Jardim Botânico, Rio de Janeiro,
RJ 22460-320, Brazil*

Received December 4, 2000; Revised July 12, 2001; Accepted October 30, 2001

Abstract. In the parallel variable distribution framework for solving optimization problems (PVD), the variables are distributed among parallel processors with each processor having the primary responsibility for updating its block of variables while allowing the remaining “secondary” variables to change in a restricted fashion along some easily computable directions. For constrained nonlinear programs convergence theory for PVD algorithms was previously available only for the case of convex feasible set. Additionally, one either had to assume that constraints are block-separable, or to use exact projected gradient directions for the change of secondary variables. In this paper, we propose two new variants of PVD for the constrained case. Without assuming convexity of constraints, but assuming block-separable structure, we show that PVD subproblems can be solved inexactly by solving their quadratic programming approximations. This extends PVD to nonconvex (separable) feasible sets, and provides a constructive practical way of solving the parallel subproblems. For inseparable constraints, but assuming convexity, we develop a PVD method based on suitable approximate projected gradient directions. The approximation criterion is based on a certain error bound result, and it is readily implementable. Using such approximate directions may be especially useful when the projection operation is computationally expensive.

Keywords: parallel optimization, variable distribution, constrained optimization, sequential quadratic programming, projected gradient

1. Introduction and motivation

We consider parallel algorithms for solving constrained optimization problems

$$\min_{x \in \mathcal{C}} f(x), \quad (1.1)$$

where \mathcal{C} is a nonempty closed set in \mathfrak{R}^n , and $f : \mathfrak{R}^n \rightarrow \mathfrak{R}$ is a continuously differentiable function. Our approach consists of partitioning the problem variables $x \in \mathfrak{R}^n$ into p blocks x_1, \dots, x_p , such that

$$x = (x_1, \dots, x_p), \quad x_l \in \mathfrak{R}^{n_l}, \quad l = 1, \dots, p, \quad \sum_{l=1}^p n_l = n,$$

*Research of the first author is supported by CNPq and FAPERJ. The second author is supported in part by CNPq Grant 300734/95-6, by PRONEX–Optimization, and by FAPERJ.

Note that if I denotes the identity matrix of appropriate dimension, the linear transformation

$$A_l^i := \begin{pmatrix} I_{n_l \times n_l} & 0_{n_l \times p-1} \\ 0_{n_l \times n_l} & D_l^i \end{pmatrix}$$

maps $(x_l, \mu_l) \in \mathfrak{R}^{n_l} \times \mathfrak{R}^{p-1}$ into $(x_l, D_l^i \mu_l) \in \mathfrak{R}^n$. In the unconstrained case, more general transformations can be used, which give rise to a fairly broad *parallel variable transformation* framework discussed in [7]. However, the theory of [7] does not appear to extend to the constrained case, which is the focus of the present paper.

We describe next the basic PVD algorithm [6, 19, 20].

Algorithm 1 (PVD). *Start with any $x^0 \in \mathcal{C}$. Choose a PVD-direction $d^0 \in \mathfrak{R}^n$. Set $i := 0$. Having x^i , check a stopping criterion. If it is not satisfied, compute x^{i+1} as follows.*

Parallelization. For each processor $l \in \{1, \dots, p\}$ compute a (possibly approximate) solution $(y_l^i, \mu_l^i) \in \mathfrak{R}^{n_l} \times \mathfrak{R}^{p-1}$ of

$$(P_l) \quad \begin{cases} \min_{x_l, \mu_l} \psi_l^i(x_l, \mu_l) := f(x_l, x_l^i + D_l^i \mu_l) \\ (x_l, x_l^i + D_l^i \mu_l) \in \mathcal{C}. \end{cases}$$

Synchronization. Compute $x^{i+1} \in \mathcal{C}$ such that

$$f(x^{i+1}) \leq \min_{l \in \{1, \dots, p\}} \psi_l^i(y_l^i, \mu_l^i).$$

Set $i := i + 1$; choose new PVD-direction d^i ; and repeat.

The idea of PVD algorithm is to balance the reduction in the number of variables for each (P_l) with allowing just enough freedom for the change of other (secondary) variables. Due to this, the parallel subproblems better approximate the original problem that has to be resolved. Note that because the secondary variables can change only along chosen fixed directions, their inclusion does not significantly increase the dimensionality of the parallel subproblems (P_l) . Indeed, the number of variables in (P_l) is $n_l + p - 1$, only $p - 1$ more than if the secondary variables were excluded. Of course, in the context of parallel computing it is reasonable to assume that p is small relative to n_l .

The synchronization step in Algorithm 1 may consist of minimizing the objective function in the affine hull, subject to feasibility, of all the points computed in parallel by the p processors. This would require solving a p -dimensional problem, which is again small compared to the original one. If (1.1) is a convex program, one can alternatively define x^{i+1} as a convex combination of the candidate points. In principle, for convergence purposes, any point with the objective function value at least as good as the smallest computed by all the processors is acceptable. As for PVD-directions, they are typically some easily computable feasible descent directions for the objective function f at the current iterate x^i , e.g., quasi-Newton or steepest descent directions in the unconstrained case.

Algorithm 1 is a rather general framework, which has to be further refined and specialized to obtain implementable/practical versions. In this respect, the two principal questions are:

- how to set up parallel subproblems; e.g., how to choose PVD-directions, and
- how to solve each parallel subproblem, including some criteria for inexact resolution.

When (1.1) is unconstrained, some of these issues have been addressed in [19], which contains improved convergence results (compared to [6]), including linear rate of convergence. These results, as well as useful generalizations such as algorithms with inexact subproblem solution and a certain degree of asynchronization, were obtained by imposing natural restrictions (of sufficient descent type) on the PVD-directions. An even more general framework for the unconstrained case was developed later in [7], where subproblems are obtained via certain nondegenerate transformations of the original variable space. These transformations can be very general, and there need not even be a distinction between primary and secondary variables. However, this approach does not seem to extend to the constrained case. It seems also that intuitively justifiable transformations do have some kind of primary-secondary variable structure. For those reasons, in the present paper we shall restrict our consideration to specific transformations of the form (1.2).

When (1.1) is a constrained optimization problem, many questions are still open, especially for a nonconvex feasible set \mathcal{C} . When \mathcal{C} is convex with block-separable structure (i.e., \mathcal{C} is a Cartesian product of closed convex sets), it was shown in [6] that every accumulation point of the PVD iterates satisfies the first-order necessary optimality conditions for problem (1.1). It was further stated that in the case of inseparable convex constraints, the PVD approach may fail. This conclusion was supported by a counter-example, which we reproduce below:

$$\begin{aligned} \min \quad & x_1^2 + x_2^2 \\ \text{s.t.} \quad & x_1 + x_2 - 2 \geq 0. \end{aligned}$$

This strongly convex quadratic program has the unique global solution at $\bar{x} = (1, 1)$. Consider any point $\hat{x} \neq \bar{x}$ such that $\hat{x}_1 + \hat{x}_2 = 2$, $\hat{x}_1 \geq 0$, $\hat{x}_2 \geq 0$, and observe that

$$\begin{aligned} \hat{x}_1 &= \arg \min_{x_1} \{x_1^2 + \hat{x}_2^2 \mid x_1 + \hat{x}_2 - 2 \geq 0\}, \\ \hat{x}_2 &= \arg \min_{x_2} \{\hat{x}_1^2 + x_2^2 \mid \hat{x}_1 + x_2 - 2 \geq 0\}. \end{aligned}$$

Therefore, if we apply Algorithm 1 using \hat{x} as its starting point and fixing the secondary variables, then this PVD variant will stay at this same nonoptimal point \hat{x} , thus failing to solve the original problem. This shows that in the constrained case, a special care should be taken in setting up the parallel subproblems.

For a general convex feasible set, it was shown in [20] that using the projected gradient direction $d(x) := x - P_{\mathcal{C}}[x - \nabla f(x)]$ for secondary variables does the job (here $P_{\mathcal{C}}[\cdot]$ stands for the orthogonal projection map onto the closed convex set \mathcal{C}). Specifically, it was established that setting $d^i := d(x^i)$ would ensure convergence of PVD methods for problems

with general (inseparable) convex constraints. Some criteria for inexact subproblem solution were also given in [20]. When \mathcal{C} is a polyhedral set, computing the projected gradient direction requires solving at every synchronization step of Algorithm 1 a single quadratic programming problem. For this, a wealth of fast and reliable algorithms is available [5, 11]. It should be noted that in the case of nonlinear constraints, the task of computing the projected gradient directions is considerably more computationally expensive. Actually, even in the affine case computing those directions exactly (or very accurately) may turn to be rather wasteful, especially when far from the solution of the original problem. Therefore, improvements are necessary.

In this paper, we propose two new versions of PVD for the nonlinearly constrained case. The first one applies to problems with block-separable *nonconvex* feasible sets. It is based on the use of *sequential quadratic programming* techniques (SQP). Our second proposal is for *inseparable* convex feasible sets. We introduce a computable approximation criterion which allows to employ *inexact* projected gradient directions. This criterion is based on an error bound result, which is of independent interest. We emphasize that its is readily implementable and preserves global convergence of PVD methods based on exact directions.

Our notation is fairly standard. The usual inner product of two vectors $x, y \in \mathfrak{R}^n$ is denoted by $\langle x, y \rangle$, and the associated norm is given by $|x|^2 = \langle x, x \rangle$. When using other norms, we shall specify them explicitly. Analogous notation will be used for subspaces of any other dimensions; for example, the reduced subspaces \mathfrak{R}^{n_l} . For the sake of simplicity, we sometimes use a compact (transposed) notation when referring to composite vectors. For instance, (x_l, μ_j) stands for the column vector $\begin{pmatrix} x_l \\ \mu_j \end{pmatrix}$. For a differentiable function $f: \mathfrak{R}^n \rightarrow \mathfrak{R}$, $\nabla f(x)$ will denote the n -dimensional column vector of partial derivatives of f at the point $x \in \mathfrak{R}^n$. For a differentiable vector-function $c: \mathfrak{R}^n \rightarrow \mathfrak{R}^m$, $c'(x)$ will denote the $m \times n$ Jacobian matrix whose rows are the transposed gradients $\nabla c_j(x)^\top$, $j = 1, \dots, m$ of the components of c . For a function h we write $h \in C_L^{1,1}(X)$ if its partial derivatives are Lipschitz-continuous on the set X with modulus $L > 0$.

2. Nonconvex separable constraints

Suppose the feasible set \mathcal{C} in (1.1) is described by a system of inequality constraints:

$$\mathcal{C} := \{x \in \mathfrak{R}^n \mid c(x) \leq 0\}, \quad (2.3)$$

where $c: \mathfrak{R}^n \rightarrow \mathfrak{R}^m$. In this section, we assume that \mathcal{C} has block-separable structure. Specifically,

$$\begin{aligned} c(x) &= (c_1(x_1), \dots, c_p(x_p)), \quad \text{where } c_l: \mathfrak{R}^{n_l} \rightarrow \mathfrak{R}^{m_l}, \\ l &= 1, \dots, p, \quad \sum_{l=1}^p m_l = m. \end{aligned} \quad (2.4)$$

Our proposal is to solve parallel subproblems (P_l) of Algorithm 1 inexactly, by making one step of the sequential quadratic programming method (SQP) [3, Ch. 13]. Since SQP methods

are local by nature, we modify the synchronization step in Algorithm 1 by introducing a suitable line-search based on the following exact penalty function [9]:

$$\theta_\sigma(x) := f(x) + \sigma \sum_{l=1}^p |c_l^+(x_l)|_1, \quad (2.5)$$

where σ is a positive parameter, and $y^+ := \max\{0, y\}$ with maximum taken componentwise. Given x^i , the l -th block of $\nabla f(x^i)$ will be denoted by $g_l^i := (I_{n_l \times n_l} \ 0_{n_l \times n_l}) \nabla f(x^i)$. Our algorithm is the following.

Algorithm 2. Start with any $x^0 \in \mathcal{C}$. Choose parameters $\underline{\sigma} > 0$ and $\tau, \beta \in (0, 1)$, and positive definite $n_l \times n_l$ matrices $M_l^0, l = 1, \dots, p$. Set $i := 0$.

Having x^i , check a stopping criterion. If it is not satisfied, compute x^{i+1} as follows.

Parallelization. For each processor $l \in \{1, \dots, p\}$ compute $(\delta_l^i, \lambda_l^i) \in \mathfrak{R}^{n_l} \times \mathfrak{R}^{m_l}$ as a KKT point of

$$(QP_l) \quad \begin{cases} \min_{\delta_l \in \mathfrak{R}^{n_l}} \langle g_l^i, \delta_l \rangle + \frac{1}{2} \langle \delta_l, M_l^i \delta_l \rangle \\ c_l(x_l^i) + c_l'(x_l^i) \delta_l \leq 0. \end{cases}$$

Synchronization. Define

$$\delta^i := (\delta_1^i, \dots, \delta_p^i) \in \mathfrak{R}^n \quad \text{and} \quad \lambda^i := (\lambda_1^i, \dots, \lambda_p^i) \in \mathfrak{R}^m.$$

Line-search. Choose $\sigma_i \geq |\lambda^i|_\infty + \underline{\sigma}$ and compute

$$\Delta^i := \langle \nabla f(x^i), \delta^i \rangle - \sigma_i \sum_{l=1}^p |c_l^+(x_l^i)|_1.$$

Using the merit function (2.5), find m_i , the smallest nonnegative integer m , such that

$$\theta_{\sigma_i}(x^i + \beta^m \delta^i) \leq \theta_{\sigma_i}(x^i) + \tau \beta^m \Delta^i. \quad (2.6)$$

Set $t_i := \beta^{m_i}$, $x^{i+1} := x^i + t_i \delta^i$, $i := i + 1$. Choose new matrices $M_l^i, l = 1, \dots, p$, and repeat.

The above proposal has to be compared to the original PVD algorithm [6]. Two remarks are in order. First, Algorithm 2 can be viewed as a PVD method with inexact solution of parallel subproblems. Indeed, “hard” nonlinear PVD subproblems (P_l) of Algorithm 1 are approximated here by “easy” quadratic programming subproblems (QP_l). And secondly, the PVD approach is extended to the case of nonconvex constraints. Note that the inclusion of forget-me-not terms does not seem to be crucial in the block-separable case (for example, such terms were not specified in the analysis of [6]). Thus, we drop here secondary variables

in (P_l) , by taking null PVD-directions. However, the use of secondary variables deserves further study for feasible sets with inseparable constraints. We also note that Algorithm 2 can be thought of as a distributed parallel implementation of SQP, where a block-diagonal matrix is chosen to generate quadratic approximations of the objective function. We remark that the contribution of Algorithm 2 is meant primarily to PVD framework rather than general SQP methods.

In Algorithm 2, we assume that subproblems (QP_l) have nonempty feasible sets for every iteration (this is guaranteed, for example, if for each l the Jacobian c'_l maps \mathfrak{N}^{n_l} onto \mathfrak{N}^{m_l} , or if c_l is convex). Alternatively, it is known that feasibility of subproblems can be forced by introducing an extra “slack” variable [1, p. 377]. It is sometimes argued that SQP methods are not convenient for solving large-scale (with n and m large) nonlinear programs when there are inequality constraints present. More precisely, it is argued that the combinatorial aspect introduced by the complementarity condition in the KKT system associated to each QP makes the subproblems resolution relatively costly. On the other hand, SQP techniques are known to be very efficient for small to medium size problems, and they are often the method of choice in that setting. In relation to Algorithm 2, it is important to note that each subproblem (QP_l) is a quadratic programming problem of relatively small dimension (n_l and m_l are presumed to be small compared to n and m). There exist a number of very fast and reliable algorithms for solving such problems (see, e.g., [5]). In Section 4, we report some numerical results for Algorithm 2, obtained via simulation on a serial computer.

For further reference, we state the optimality conditions for (QP_l) : the pair $(\delta_l^i, \lambda_l^i) \in \mathfrak{N}^{n_l} \times \mathfrak{N}^{m_l}$ solves the KKT system

$$g_l^i + M_l^i \delta_l^i + (c'_l(x_l^i))^\top \lambda_l^i = 0, \quad (2.7a)$$

$$c_l(x_l^i) + c'_l(x_l^i) \delta_l^i \leq 0, \quad (2.7b)$$

$$\lambda_l^i \geq 0 \quad \text{and} \quad (\lambda_l^i, c_l(x_l^i) + c'_l(x_l^i) \delta_l^i) = 0. \quad (2.7c)$$

Denoting by $\theta'_\sigma(x; d)$ the usual directional derivative of the merit function θ_σ at $x \in \mathfrak{N}^n$ in the direction $d \in \mathfrak{N}^n$, we next state that δ^i is a descent direction for this function at the point x^i . This in turn will imply that the line-search procedure in Algorithm 2 is well-defined.

Lemma 1. *If $f, c \in C_L^{1,1}(\mathfrak{N}^n)$, then*

$$\theta'_{\sigma_i}(x^i; \delta^i) \leq \Delta^i \leq -\sum_{l=1}^p \langle \delta_l^i, M_l^i \delta_l^i \rangle - \sigma \sum_{l=1}^p |c_l^+(x_l^i)|_1, \quad (2.8)$$

where $x^i, \delta^i, \Delta^i, M_l^i$ are defined in Algorithm 2.

Proof: Defining the index-sets $I_+(x^i) := \{j \mid c_j(x^i) > 0\}$, $I_=(x^i) := \{j \mid c_j(x^i) = 0\}$, we have that (e.g., see [9, p. 301])

$$\theta'_{\sigma_i}(x^i; \delta^i) = \langle \nabla f(x^i), \delta^i \rangle + \sigma_i \sum_{j \in I_+(x^i)} \langle \nabla c_j(x^i), \delta^i \rangle + \sigma_i \sum_{j \in I_=(x^i)} (\langle \nabla c_j(x^i), \delta^i \rangle)^+.$$

Using (2.7b) componentwise, we have

$$\langle \nabla c_j(x^i), \delta^i \rangle \leq -c_j(x^i) \begin{cases} < 0 & \text{if } j \in I_+(x^i) \\ = 0 & \text{if } j \in I_=(x^i) \end{cases}$$

implying that

$$\begin{aligned} \langle \nabla c_j(x^i), \delta^i \rangle &\leq -|c_j^+(x^i)|, & j \in I_+(x^i), \\ (\langle \nabla c_j(x^i), \delta^i \rangle)^+ &\leq -|c_j^+(x^i)|, & j \in I_=(x^i). \end{aligned}$$

Hence,

$$\theta'_{\sigma_i}(x^i; \delta^i) \leq \langle \nabla f(x^i), \delta^i \rangle - \sigma_i \sum_{l=1}^p |c_l^+(x^i)|_1 = \Delta^i. \quad (2.9)$$

To obtain the right-most inequality in (2.8) we proceed as follows:

$$\begin{aligned} \langle \nabla f(x^i), \delta^i \rangle &= \sum_{l=1}^p \langle g_l^i, \delta_l^i \rangle \\ &= -\sum_{l=1}^p \langle M_l^i \delta_l^i + (c_l'(x^i))^\top \lambda_l^i, \delta_l^i \rangle \\ &= -\sum_{l=1}^p \langle M_l^i \delta_l^i, \delta_l^i \rangle - \sum_{l=1}^p \langle \lambda_l^i, c_l'(x^i) \delta_l^i \rangle \\ &= -\sum_{l=1}^p \langle M_l^i \delta_l^i, \delta_l^i \rangle + \sum_{l=1}^p \langle \lambda_l^i, c_l(x^i) \rangle \\ &\leq -\sum_{l=1}^p \langle M_l^i \delta_l^i, \delta_l^i \rangle + \sum_{l=1}^p \langle \lambda_l^i, c_l^+(x^i) \rangle, \end{aligned}$$

where the second equality is by (2.7a), the fourth is by (2.7c), and the inequality follows from (2.7c) and the fact that $c_l^+(x^i) \geq c_l(x^i)$.

Combining the latter relation with (2.9), we now have that

$$\begin{aligned} \Delta^i &\leq -\sum_{l=1}^p \langle M_l^i \delta_l^i, \delta_l^i \rangle + \sum_{l=1}^p \langle \lambda_l^i, c_l^+(x^i) \rangle - \sigma_i \sum_{l=1}^p |c_l^+(x^i)|_1 \\ &\leq -\sum_{l=1}^p \langle M_l^i \delta_l^i, \delta_l^i \rangle + \sum_{l=1}^p (|\lambda_l^i|_\infty - \sigma_i) |c_l^+(x^i)|_1 \\ &\leq -\sum_{l=1}^p \langle M_l^i \delta_l^i, \delta_l^i \rangle - \underline{\sigma} \sum_{l=1}^p |c_l^+(x^i)|_1, \end{aligned}$$

where the second inequality is by the Cauchy-Schwarz inequality, and the last is by the choice of σ_i . This completes the proof. \square

For Algorithm 2 to be globally convergent, we assume that the penalization parameter σ_i is kept bounded above, which essentially means that the multipliers λ_l^i stay bounded. From the practical point of view, the latter assumption is natural. In what follows, we show convergence of Algorithm 2 to Karush-Kuhn-Tucker (KKT) points of (1.1), i.e., pairs $(\bar{x}, \bar{\lambda}) \in \mathfrak{R}^n \times \mathfrak{R}^m$ satisfying

$$\begin{aligned} \nabla f(\bar{x}) + c'(\bar{x})^\top \bar{\lambda} &= 0, \\ c(\bar{x}) &\leq 0, \\ \bar{\lambda} &\geq 0 \quad \text{and} \quad \langle \bar{\lambda}, c(\bar{x}) \rangle = 0. \end{aligned} \tag{2.10}$$

Theorem 1. *Suppose that $f, c \in C_L^{1,1}(\mathfrak{R}^n)$, and let the feasible set \mathcal{C} have block-separable structure given by (2.4). Let $\{(x^i, \lambda^i)\}$ be a sequence generated by Algorithm 2. Assume there exists an iteration index i_0 such that $\sigma_i = \bar{\sigma}$ for all $i \geq i_0$ and the matrices M_l^i are uniformly positive definite and bounded for all $l = 1, \dots, p$ and all iteration indices i . Then either $\theta_{\sigma_i}(x^i) \rightarrow -\infty$, or every accumulation point $(\bar{x}, \bar{\lambda})$ of the sequence $\{(x^i, \lambda^i)\}$ is a KKT point of the problem, i.e., it satisfies (2.10).*

Proof: If for some iteration index i it happens that

$$\delta_l^i = 0 \quad \text{and} \quad c_l^+(x_l^i) = 0, \quad l = 1, \dots, p, \tag{2.11}$$

then (2.7a)–(2.7c) reduce to

$$\begin{aligned} g_l^i + (c_l'(x_l^i))^\top \lambda_l^i &= 0, \\ c_l(x_l^i) &\leq 0, \\ \lambda_l^i &\geq 0 \quad \text{and} \quad \langle \lambda_l^i, c_l(x_l^i) \rangle = 0. \end{aligned}$$

Now taking into account separability of constraints, it follows that (x^i, λ^i) satisfies the KKT system (2.10).

Suppose now that (2.11) does not hold for any i . By (2.8) in Lemma 1, δ^i is then a direction of descent for the merit function θ_{σ_i} at x^i . By standard argument, the line-search procedure is well-defined and terminates finitely with some stepsize $t_i > 0$. The entire method is then well-defined and generates an infinite sequence of iterates.

We prove first that the sequence of stepsizes $\{t_i\}$ is bounded away from 0. Take any $t \in (0, 1]$. Since $f \in C_L^{1,1}(\mathfrak{R}^n)$, by [1, Proposition A.24] we have that

$$f(x^i + t\delta^i) \leq f(x^i) + t\langle \nabla f(x^i), \delta^i \rangle + \frac{Lt^2}{2}|\delta^i|^2. \tag{2.12}$$

Similarly, since $c \in C_L^{1,1}(\mathfrak{R}^n)$, using the equivalence of the norms in the finite-dimensional setting, for some $R_1 > 0$ we have that

$$\begin{aligned} R_1 t^2 |\delta_i^i|^2 &\geq |c_l(x_l^i + t\delta_l^i) - c_l(x_l^i) - t c_l'(x_l^i) \delta_l^i|_1 \\ &= |c_l(x_l^i + t\delta_l^i) - t(c_l(x_l^i) + c_l'(x_l^i) \delta_l^i) - (1-t)c_l(x_l^i)|_1 \\ &\geq |(c_l(x_l^i + t\delta_l^i) - t(c_l(x_l^i) + c_l'(x_l^i) \delta_l^i)) - (1-t)c_l(x_l^i)|_1 \\ &\geq |c_l^+(x_l^i + t\delta_l^i)|_1 - |(t(c_l(x_l^i) + c_l'(x_l^i) \delta_l^i) + (1-t)c_l(x_l^i))|_1. \end{aligned}$$

where the equality is obtained by adding and subtracting $t c_l(x_l^i)$, the second inequality follows from $|a|_1 \geq |a^+|_1$, and $|(a-b)^+|_1 \geq |a^+|_1 - |b^+|_1$ is used in the last inequality.

Re-writing the above relation, we obtain

$$\begin{aligned} |c_l^+(x_l^i + t\delta_l^i)|_1 &\leq |(t(c_l(x_l^i) + c_l'(x_l^i) \delta_l^i) + (1-t)c_l(x_l^i))|_1 + R_1 t^2 |\delta_l^i|^2 \\ &\leq t|(c_l(x_l^i) + c_l'(x_l^i) \delta_l^i)|_1 + (1-t)|c_l^+(x_l^i)|_1 + R_1 t^2 |\delta_l^i|^2 \\ &= (1-t)|c_l^+(x_l^i)|_1 + R_1 t^2 |\delta_l^i|^2, \end{aligned}$$

where the second inequality is by the convexity of $|\cdot|_1$, and the equality is by (2.7b).

Together with (2.12), the last inequality yields

$$\begin{aligned} \theta_{\sigma_i}(x^i + t\delta^i) &= f(x^i + t\delta^i) + \sigma_i \sum_{l=1}^p |c_l^+(x_l^i + t\delta_l^i)|_1 \\ &\leq f(x^i) + t \langle \nabla f(x^i), \delta^i \rangle + \sigma_i (1-t) \sum_{l=1}^p |c_l^+(x_l^i)|_1 + R_2 t^2 |\delta^i|^2 \\ &= \theta_{\sigma_i}(x^i) + t\Delta^i + R_2 t^2 |\delta^i|^2, \end{aligned}$$

where $R_2 > 0$ is a fixed constant depending on $R_1, L, \bar{\sigma}$. By a direct comparison of the latter relation with (2.6), we conclude that (2.6) is guaranteed to be satisfied once m is large enough so that $t = \beta^m$ falls within the set of t satisfying $t\Delta^i + R_2 t^2 |\delta^i|^2 \leq \tau t\Delta^i$. In particular, since the line-search procedure did not accept the stepsize $t = \beta^{m_i-1}$, it follows that

$$\text{either } m_i = 0 \text{ or } t = \beta^{m_i-1} > \frac{(\tau-1)\Delta^i}{R_2 |\delta^i|^2}. \quad (2.13)$$

By (2.8) in Lemma 1 and the uniform positive definiteness of matrices M_l^i , there exists $R_3 > 0$ such that

$$-\Delta^i \geq \sum_{l=1}^p \langle M_l^i \delta_l^i, \delta_l^i \rangle \geq \sum_{l=1}^p R_3 |\delta_l^i|^2 = R_3 |\delta^i|^2,$$

and using (2.13), we conclude that

$$t_i = \beta^{m_i} \geq \bar{t} := \min\{1; \beta(1 - \tau)R_3/R_2\} > 0.$$

By the assumption that $\sigma_i = \bar{\sigma}$ for all $i \geq i_0$, from (2.6) it follows that the sequence $\{\theta_{\bar{\sigma}}(x^i)\}$ is nonincreasing. Hence, it is either unbounded below, or it converges. In the latter case, (2.6) implies that $\tau t_i \Delta^i \rightarrow 0$, and since $t_i \geq \bar{t} > 0$, it holds that

$$\Delta^i \rightarrow 0.$$

By the definition of Δ^i and the uniform positive definiteness of matrices M_l^i , we conclude that

$$\delta_l^i \rightarrow 0, \quad c_l^+(x_l^i) \rightarrow 0, \quad l = 1, \dots, p.$$

Now, passing onto the limit in (2.7a)–(2.7c) as $i \rightarrow \infty$, and taking into account the boundedness of the matrices M_l^i , we obtain the assertions of the theorem. \square

3. Convex inseparable constraints

Suppose now that the feasible set \mathcal{C} is defined by a system of convex inequalities, i.e., $c : \mathfrak{R}^n \rightarrow \mathfrak{R}^m$ in (2.3) has convex components $c_j(\cdot)$, $j = 1, \dots, m$. We emphasize that in this section we do not assume separability of constraints. In this setting, it appears that the only currently known way to ensure convergence of the PVD algorithm is to use the projected gradient directions for the change of secondary variables [20]. Omitting the iteration indices i , if x is the current iterate, to compute these directions one has to solve a subproblem of the following structure:

$$\min_{c(z) \leq 0} \frac{1}{2} |z - (x - \nabla f(x))|^2. \quad (3.14)$$

This is done by some iterative algorithm. As already discussed in Section 1, solving this problem in the general nonlinear case can be quite costly. Moreover, when far from a solution of the original problem, exact (or even very accurate) projection directions are perhaps unnecessary. This suggests developing a stopping rule for solving (3.14) or, equivalently, an approximation criterion for projection directions to be used in the PVD scheme. For algorithmic purposes, it is important to make this criterion constructive and implementable.

Assuming some constraint qualification condition [13], we have that \bar{z} solves (3.14), i.e.,

$$\bar{z} = P_{\mathcal{C}}[x - \nabla f(x)]$$

if, and only if, the pair $(\bar{z}, \bar{u}) \in \mathfrak{R}^n \times \mathfrak{R}^m$ satisfies the KKT system

$$\begin{aligned} \nabla_z L(\bar{z}, \bar{u}) &= \bar{z} - x + \nabla f(x) + c'(\bar{z})^\top \bar{u} = 0, \\ c(\bar{z}) &\leq 0, \\ \bar{u} &\geq 0 \quad \text{and} \quad \langle \bar{u}, c(\bar{z}) \rangle = 0, \end{aligned} \tag{3.15}$$

where

$$L(z, u) = \frac{1}{2} |z - (x - \nabla f(x))|^2 + \langle u, c(z) \rangle \tag{3.16}$$

is the standard Lagrangian for problem (3.14).

Suppose $z \in \mathcal{C}$ and $u \in \mathfrak{R}_+^m$ is some current approximation to a primal-dual optimal solution of (3.14), generated by an iterative algorithm applied to solve this problem. Lemma 2 below establishes an *error bound* for the distance from z to $P_{\mathcal{C}}[x - \nabla f(x)]$ in terms of violations of the KKT conditions (3.15) by the pair (z, u) . A nice feature of this estimate is that unlike some (perhaps, most) error bounds results in the literature (see [18] for a survey), it does not involve any expressions which are not readily computable or any other quantities which are not *observable*. Furthermore, this error bound holds globally, i.e., not just in some neighbourhood of the solution point. Thus it can be easily employed for algorithmic purposes.

Lemma 2 is related to error bounds for strongly convex programs obtained in [15]. However, Theorem 2.2 in [15] involves certain constants which are in general not computable, while Corollary 2.4 in [15] assumes not only that z is primal feasible, but also that (z, u) is *dual* feasible, which here means that $u \geq 0$, $\nabla_z L(z, u) = 0$. In Lemma 2 we only assume that z is primal feasible and that the approximate multiplier u is nonnegative. Our assumptions are not only weaker but they also appear to be more suitable in an algorithmic framework, as they will be satisfied at each iteration of many standard optimization methods. Dual feasibility, in contrast, is unlikely to be satisfied along iterations of typical algorithms, except in the limit.

Lemma 2. *Let $c : \mathfrak{R}^n \rightarrow \mathfrak{R}^m$ be convex and differentiable, and suppose that the set \mathcal{C} given by (2.3) satisfies some constraint qualification.*

Then for any $z \in \mathcal{C}$ and $u \in \mathfrak{R}_+^m$, it holds that

$$|z - P_{\mathcal{C}}[x - \nabla f(x)]| \leq \varepsilon(z, u),$$

where

$$\varepsilon(z, u) := \frac{1}{2} |\nabla_z L(z, u)| + \frac{1}{2} \sqrt{|\nabla_z L(z, u)|^2 - 4\langle u, c(z) \rangle}. \tag{3.17}$$

Proof: Let $\bar{z} = P_{\mathcal{C}}[x - \nabla f(x)]$ and \bar{u} be the associated multiplier, so that the pair (\bar{z}, \bar{u}) satisfies (3.15). Denote

$$\gamma := \langle \nabla_z L(z, u) - \nabla_z L(\bar{z}, \bar{u}), z - \bar{z} \rangle - \langle c(z) - c(\bar{z}), u - \bar{u} \rangle,$$

where $L(\cdot, \cdot)$ is defined by (3.16). Take any $z \in \mathcal{C}$ and $u \in \mathfrak{N}_+^m$. Denoting further $y = x - \nabla f(x)$, we have that

$$\begin{aligned} \gamma &= \langle z - y + c'(z)^\top u - (\bar{z} - y + c'(\bar{z})^\top \bar{u}), z - \bar{z} \rangle - \langle c(z) - c(\bar{z}), u - \bar{u} \rangle \\ &= |z - \bar{z}|^2 + \langle u, c(\bar{z}) - c(z) - c'(z)(\bar{z} - z) \rangle + \langle \bar{u}, c(z) - c(\bar{z}) - c'(\bar{z})(z - \bar{z}) \rangle \\ &\geq |z - \bar{z}|^2, \end{aligned} \quad (3.18)$$

where the inequality follows from $u \geq 0$, $\bar{u} \geq 0$ and the facts that, by the convexity of $c(\cdot)$, $c(\bar{z}) - c(z) - c'(z)(\bar{z} - z) \geq 0$ and $c(z) - c(\bar{z}) - c'(\bar{z})(z - \bar{z}) \geq 0$.

On the other hand,

$$\begin{aligned} \gamma &= \langle \nabla_z L(z, u), z - \bar{z} \rangle - \langle c(z), u \rangle + \langle \bar{u}, c(z) \rangle + \langle u, c(\bar{z}) \rangle \\ &\leq |\nabla_z L(z, u)| |z - \bar{z}| - \langle u, c(z) \rangle, \end{aligned} \quad (3.19)$$

where the equality follows from the KKT conditions (3.15), and the inequality follows from the facts that $\bar{u} \geq 0$, $u \geq 0$ and $c(z) \leq 0$, $c(\bar{z}) \leq 0$, and the Cauchy-Schwarz inequality. Denoting $t := |z - \bar{z}|$, $\alpha := |\nabla_z L(z, u)| \geq 0$ and $\beta := -\langle u, c(z) \rangle \geq 0$, and combining (3.18) with (3.19), we obtain the following quadratic inequality in t :

$$t^2 - \alpha t - \beta \leq 0.$$

Resolving this inequality, we obtain that

$$t \leq \frac{1}{2}(\alpha + \sqrt{\alpha^2 + 4\beta}).$$

Recalling definitions of the quantities involved, we conclude that

$$|z - \bar{z}| \leq \frac{1}{2}|\nabla_z L(z, u)| + \frac{1}{2}\sqrt{|\nabla_z L(z, u)|^2 - 4\langle u, c(z) \rangle} = \varepsilon(z, u). \quad \square$$

We propose the following PVD algorithm based approximations of the projected gradient directions.

Algorithm 3. Choose parameters $\sigma_1 \in (0, 1)$ and $\sigma_2 \in (0, (1 - \sigma_1)^2)$. Start with any $x^0 \in \mathcal{C}$. Set $i := 0$.

Having x^i , check a stopping criterion. If it is not satisfied, proceed as follows.

PVD-direction choice. Compute $z^i \approx P_{\mathcal{C}}[x^i - \nabla f(x^i)]$ such that $z^i \in \mathcal{C}$ and the associated approximate Lagrange multiplier $u^i \in \mathfrak{N}_+^m$ for Problem (3.14) satisfy

$$\varepsilon(z^i, u^i) \leq \min \left\{ \sigma_1, \frac{\sigma_2 |z^i - x^i|}{|\nabla f(x^i)|} \right\} |z^i - x^i|, \quad (3.20)$$

where $\varepsilon(z^i, u^i)$ is given by (3.17). Set $d^i := x^i - z^i$.

Compute x^{i+1} as follows.

Parallelization. For each processor $l \in \{1, \dots, p\}$ compute a solution $(y_l^i, \mu_l^i) \in \mathfrak{R}^{n_l} \times \mathfrak{R}^{p-1}$ of (P_l) as defined in Algorithm 1.

Synchronization. Compute $x^{i+1} \in \mathcal{C}$ such that

$$f(x^{i+1}) \leq \min_{l \in \{1, \dots, p\}} \psi_l^i(y_l^i, \mu_l^i).$$

Set $i := i + 1$; and repeat.

Note that Algorithm 3.1 is a general-purpose method for problems with no special structure. The example presented in the Introduction shows that for such problems computing a meaningful direction for secondary variables is indispensable. Compared to *any* alternative, computing an approximate projected gradient direction seems to be quite favorable in terms of cost, and perhaps even the best one can do. Regarding the tolerance criterion (3.20), we note that if z^i is the exact projection point then $\varepsilon(z^i, u^i) = 0$, while $z^i - x^i \neq 0$ (if $z^i = x^i$ then x^i satisfies the first-order necessary optimality condition $x^i = P_{\mathcal{C}}[x^i - \nabla f(x^i)]$). From this, it is easy to see that the projection problem (3.14) always has inexact solutions satisfying (3.20). Hence, the method is well-defined.

As for the convergence properties of Algorithm 3, our result is the following.

Theorem 2. *Let $c : \mathfrak{R}^n \rightarrow \mathfrak{R}^m$ be convex and differentiable, and $f \in C_L^{1,1}(\mathcal{C})$. Suppose $\{x^i\}$ is a sequence generated by Algorithm 3. Then either f is unbounded from below on \mathcal{C} , or the sequence $\{f(x^i)\}$ converges, and every accumulation point of the sequence $\{x^i\}$ satisfies the first-order necessary optimality condition.*

Proof: Consider any iteration $i = 0, 1, \dots$, any processor $l \in \{1, \dots, p\}$, and the point

$$(x_l^i - \alpha d_l^i, -\alpha e_l^i) \in \mathfrak{R}^{n_l+p-1}, \quad 0 < \alpha < \min \left\{ 1, \frac{2}{L}((1 - \sigma_1)^2 - \sigma_2) \right\},$$

where $e_l^i = (1, \dots, 1) \in \mathfrak{R}^{p-1}$. We first show that this point is feasible for the corresponding subproblem (P_l) of minimizing the function $\psi_l^i(x_l, \mu_l^i) := f(x_l, x_l^i + D_l^i \mu_l^i)$. Indeed, using the notation (1.2),

$$\begin{aligned} (x_l^i - \alpha d_l^i, x_l^i - \alpha D_l^i e_l^i) &= (x_l^i - \alpha d_l^i, x_l^i - \alpha d_l^i) \\ &= x^i - \alpha d^i \\ &= (1 - \alpha)x^i + \alpha z^i \in \mathcal{C}, \end{aligned}$$

where the first equality follows from the structure of D_l^i , and the inclusion is by the facts that $x^i, z^i \in \mathcal{C}$, $\alpha \in (0, 1)$, and the convexity of the set \mathcal{C} . As a result, we have that

$$\begin{aligned}
f(x^i) - f(y_l^i, x_l^i + D_l^i \mu_l^i) &= \psi_l^i(x_l^i, 0) - \psi_l^i(y_l^i, \mu_l^i) \\
&\geq \psi_l^i(x_l^i, 0) - \psi_l^i(x_l^i - \alpha d_l^i, -\alpha e_l^i) \\
&= f(x^i) - f(x_l^i - \alpha d_l^i, x_l^i - \alpha D_l^i e_l^i) \\
&= f(x^i) - f(x^i - \alpha d^i) \\
&\geq \alpha \langle \nabla f(x^i), d^i \rangle - \frac{L}{2} \alpha^2 |d^i|^2,
\end{aligned} \tag{3.21}$$

where the last inequality follows from $f \in C_L^{1,1}(\mathcal{C})$ (e.g., [1, Proposition A.24]).

By the construction of the algorithm and Lemma 2,

$$|z^i - P_C[x^i - \nabla f(x^i)]| \leq \varepsilon(z^i, u^i).$$

Hence, there exists some $\xi^i \in \mathfrak{R}^n$ such that

$$z^i + \xi^i = P_C[x^i - \nabla f(x^i)], \quad |\xi^i| \leq \varepsilon(z^i, u^i). \tag{3.22}$$

Define the (continuous) residual function

$$R(x) := x - P_C[x - \nabla f(x)].$$

With this notation,

$$x^i - z^i - \xi^i = d^i - \xi^i = R(x^i).$$

By properties of the projection operator (e.g., see [1, Proposition B.11]), since $x^i \in \mathcal{C}$ we have that

$$\begin{aligned}
0 &\geq \langle x^i - \nabla f(x^i) - P_C[x^i - \nabla f(x^i)], x^i - P_C[x^i - \nabla f(x^i)] \rangle \\
&= \langle x^i - \nabla f(x^i) - z^i - \xi^i, x^i - z^i - \xi^i \rangle \\
&= -\langle \nabla f(x^i), d^i \rangle + \langle \nabla f(x^i), \xi^i \rangle + |d^i - \xi^i|^2.
\end{aligned}$$

Hence,

$$\langle \nabla f(x^i), d^i \rangle \geq |R(x^i)|^2 - |\nabla f(x^i)| |\xi^i|.$$

Combining the latter relation with (3.21) and (3.22), and taking into account the synchronization step of Algorithm 3, we further obtain

$$f(x^i) - f(x^{i+1}) \geq \alpha \left(|R(x^i)|^2 - \varepsilon(z^i, u^i) |\nabla f(x^i)| - \frac{L}{2} \alpha |d^i|^2 \right). \tag{3.23}$$

Observe that

$$|d^i| \leq |R(x^i)| + |\xi^i| \leq |R(x^i)| + \varepsilon(z^i, u^i) \leq |R(x^i)| + \sigma_1 |d^i|,$$

where the first relation is by the Cauchy-Schwarz inequality, the second follows from (3.22), and the last from (3.20). Hence,

$$|R(x^i)| \geq (1 - \sigma_1) |d^i|.$$

Using (3.23), the latter inequality, and again (3.20), we have that

$$f(x^i) - f(x^{i+1}) \geq \alpha \left((1 - \sigma_1)^2 - \sigma_2 - \frac{L\alpha}{2} \right) |d^i|^2. \quad (3.24)$$

By the choice of α and assumptions on parameters σ_1 and σ_2 , the sequence $\{f(x^i)\}$ is nonincreasing. If $f(\cdot)$ is bounded below on \mathcal{C} , then $\{f(x^i)\}$ is bounded below, and hence it converges. In the latter case, $\{f(x^i) - f(x^{i+1})\} \rightarrow 0$ and therefore $\{d^i\} \rightarrow 0$, by (3.24). On the other hand,

$$|d^i| \geq |R(x^i)| - |\xi^i| \geq |R(x^i)| - \varepsilon(z^i, u^i) \geq |R(x^i)| - \sigma_1 |d^i|,$$

so that

$$|R(x^i)| \leq (1 + \sigma_1) |d^i|.$$

We conclude that $\{R(x^i)\}$ also tends to zero. By the continuity of $R(\cdot)$, it then follows that for every accumulation point \bar{x} of the sequence $\{x^i\}$, $R(\bar{x}) = 0$. It is well known that the latter is equivalent to the minimum principle necessary optimality condition $\bar{x} \in \mathcal{C}$, $\langle \nabla f(\bar{x}), x - \bar{x} \rangle \geq 0$ for all $x \in \mathcal{C}$. \square

4. Preliminary numerical experience

To get some insight into computational properties of our approach in Section 2, we considered test problems taken from the study of Sphere Packing Problems [4]. In particular, the problems chosen are the same as the ones used in [16]. Not having access to a parallel computer, we have carried out a simulation, i.e., the subproblems are solved serially on a serial machine. Even though this is admittedly a rather crude experiment, it nevertheless gives some idea on what one might expect in actual parallel implementation.

Given p spheres in \mathfrak{R}^v with $v, p \in \{1, 2, \dots\}$, so that $x \in \mathfrak{R}^{vp}$, the problems are as follows.

Problem 1.

$$\begin{aligned} \min \quad & \sum_{i=1}^{p-1} \sum_{j=i+1}^p \sum_{t=1}^v x_{(i-1)v+t} x_{(j-1)v+t} \\ \text{s.t.} \quad & \sum_{t=1}^v x_{(i-1)v+t}^2 - 1 = 0, \quad i = 1, \dots, p. \end{aligned}$$

Problem 2. Given an integer $\mu \geq 1$,

$$\begin{aligned} \min \quad & \sum_{i=1}^{p-1} \sum_{j=i+1}^p \frac{1}{\left(\sum_{t=1}^v (x_{(i-1)v+t} - x_{(j-1)v+t})^2 + 1 \right)^\mu} \\ \text{s.t.} \quad & \sum_{t=1}^v x_{(i-1)v+t}^2 - 1 = 0, \quad i = 1, \dots, p. \end{aligned}$$

Problem 3.

$$\begin{aligned} \min \quad & \sum_{i=1}^{p-1} \sum_{j=i+1}^p \left(\left(\sum_{t=1}^v (x_{(i-1)v+t} - x_{(j-1)v+t})^2 \right)^{-6} \right) \\ & - 2 \sum_{i=1}^{p-1} \sum_{j=i+1}^p \left(\left(\sum_{t=1}^v (x_{(i-1)v+t} - x_{(j-1)v+t})^2 \right)^{-3} \right) \\ \text{s.t.} \quad & \sum_{t=1}^v x_{(i-1)v+t}^2 - 1 = 0, \quad i = 1, \dots, p. \end{aligned}$$

We have implemented in Matlab our Algorithm 2, the serial SQP method, and the general PVD Algorithm 1. All codes were run under Matlab version 6.0.0.88 (Release 12) on a Sun UltraSPARCstation. Details of the implementation are as follows. Algorithm 1 is implemented using the function `constr.m` from the Matlab Optimization Toolbox for solving the subproblems (P_l). Algorithm 2 and the serial SQP method are quite sophisticated quasi-Newton implementations with line-search based on the merit function (2.5). Each (QP_l) is solved by the Null-Space Method [8]. The penalty parameter σ_i is updated using a modification of the Mayne and Polak rule [17] that allows σ_i to decrease, if warranted. The stepsize t_i is computed with an Armijo rule that uses safeguarded quadratic interpolation. In order to prevent the Maratos' effect (i.e., to ensure that the unit stepsize is asymptotically accepted if possible), a second-order correction is added to the search direction δ^i when necessary [3, Ch. 13.4]. Finally, the quasi-Newton matrices are updated by the BFGS formula, using also the Powell correction and the Oren-Luenberger scaling. The methods stop when the relative error, measured by the sum of the norm of the reduced gradient and the norm of the constraints, is less than 10^{-5} .

We first compare our Algorithm 2 with the general PVD Algorithm 1. Since this comparison turned out rather obvious and one-sided (which is certainly not surprising), we do

Table 1. Results for Algorithms 1 and 2 on Problem 1.

Name	p	Algorithm 1		Algorithm 2	
		Iteration	Time	Iteration	Time
Problem 1	$p = 2$	1	.81	23	1.25
$v = 3$	$p = 4$	9	19.1	13	2.35
	$p = 8$	5	153.8	10	10.8
	$p = 16$	5	287.6	9	66.5

Starting points are feasible, generated randomly. Times are measured (in seconds) using the intrinsic Matlab function `cputime`.

not report exhaustive testing for this part. In Table 1 we report the number of iterations and the running times for Algorithms 1 and 2 on Problem 1 (results for other problems follow the same trend and do not yield any further insight). Note that the running times reported are serial, without any regard to the parallel nature of the algorithms. Because the iterative structure of the two algorithms is the same, this seems to be a meaningful and fair comparison. It is already clear from intuition that solving exactly the general nonlinear subproblems (P_i) in Algorithm 1 is very costly, and is unlikely to yield a competitive algorithm. This was easily confirmed by our experience. Of course, one could use heuristic considerations to (somehow) adjust dynamically the tolerance in solving the subproblems (in our experiment, all subproblems are solved to within the tolerance of 10^{-5} , same as the original problem). However, there is no convergence analysis to support such a strategy within Algorithm 1 in the constrained case (strictly speaking, there is no even a convergence proof for Algorithm 1 in the case under consideration, since the feasible set is nonconvex!). As discussed above, one of the motivations for our Algorithm 2 is precisely to provide a constructive implementable way of solving subproblems (P_i) inexactly, by solving their quadratic approximations. The results in Table 1 confirm that the proposed approach certainly makes sense.

Our next set of experiments concerns with the comparison between Algorithm 2 and serial SQP. It is not easy to come up with a meaningful comparison of a serial method and a parallel method on a serial machine. For example, the overall running time of a “parallel” method obtained on a serial machine is considered notoriously unreliable to predict time that would be required by its parallel implementation (i.e., just dividing the time by the number of “processors” is not a good measure). We therefore focus on indicators which we feel should be still meaningful for the actual parallel implementation, as discussed below. To evaluate the gain in computing the search directions, we report the total time spent solving quadratic programming subproblems within the SQP and within the serial implementation of Algorithm 2. Since no communication between the processors would be needed within this phase, the “expected” time for solving QPs by the parallel implementation of Algorithm 2 can indeed be obtained dividing by the number of processors. It is somewhat surprising that even for smaller problems, see Table 2, the *serial* time for solving QPs in Algorithm 2 is already considerably smaller than in standard SQP, with the difference becoming drastic for larger problems. If we approximate the “speedup” efficiency for computing directions

Table 2. Results for Algorithm 2 and the SQP method on sphere packing problems.

Problem ν, μ	p	SQP			Algorithm 2			
		QP time	Iteration no.	Calls to oracle	QP time (serial)	"speedup" %	Iteration no.	Calls to oracle
Problem 1 $\nu = 4$	2	.43	11	12	.31	143	14	30
	4	8.81	10	13	5	175	31	36
	8	59.4	9	12	4.48	1323	13	20
	16	179	16	38	2.69	6656	16	23
	32	766	9	12	3.07	24951	9	10
$\nu = 10$	2	7.35	11	13	3.03	242	11	13
	4	42.4	9	12	5.5	771	10	13
	8	314	9	13	12.1	2603	11	15
	16	2167	8	13	17.6	12251	8	13
Problem 2 $\mu = 1, \nu = 4$	2	.45	12	13	.55	81	26	28
	4	13.8	15	40	2.52	549	15	22
	8	23	15	19	2.16	1056	26	27
	16	3070	65	1030	21.4	14293	32	740
$\mu = 2, \nu = 4$	32	29101	341	594	115.2	25263	341	471
	2	2.57	18	19	1.4	183	17	30
	4	17.2	20	33	3.04	570	19	34
	8	91.5	15	21	6.96	1319	21	22
	16	1385	31	59	18.4	7511	29	68
$\mu = 4, \nu = 4$	32	21865	65	76	186.2	11736	147	172
	2	2.06	50	137	.63	316	28	59
	4	95	104	265	12.8	748	76	162
	8	224	34	35	15.2	1488	44	64
	16	1807	159	215	48.96	3687	288	319
32	8077	95	155	47.7	16827	139	203	
Problem 3 $\nu = 3$	2	.91	10	11	.95	96	16	17
	4	28.6	54	89	6.95	411	58	132
	8	209	59	123	12.7	1651	53	109
	16	1522	141	223	112	1359	181	257
	32	35679	187	445	374	9551	201	511
$\nu = 4$	2	1.38	35	49	.725	194	34	66
	4	7.29	33	56	1.23	592	29	55
	8	134	88	132	12.9	1044	154	235
	16	2222	198	407	46.3	4802	274	599
	32	23289	267	630	100.5	23173	300	752

Starting points are generated randomly, with coordinates in $[-10, 10]$. Times of solving QPs are measured (in seconds) using the intrinsic Matlab function `cputime`.

in parallel by

$$\frac{\text{time spent solving QPs in SQP}}{\text{time spent solving QPs in serial Algorithm 2}} * 100,$$

then these efficiencies vary from acceptable (around 80%) to very high (over 10000%), and grow fast with the size of the problem. This confirms our motivation as discussed in Section 2, i.e., smaller QPs are significantly easier to solve. Obviously, the possible price to pay in Algorithm 2 is “deterioration” of directions compared to a good implementation of the full SQP algorithm. This issue is addressed by reporting the numbers of iterations and calls to the oracle evaluating the function and derivatives values, see Table 2. We note that the numbers of iterations and function/derivatives evaluations are usually slightly higher for Algorithm 2, but not always. Overall, those numbers for the two methods are quite similar. Given the impressive gain that Algorithm 2 exhibits in solving QP subproblems, this indicates that the parallel implementation of this algorithm should indeed be efficient, at least when computing the function and derivatives is cheap relative to solving QPs. In particular, this should be the case for large-scale problems where the functions and their derivatives are given explicitly.

5. Concluding remarks

Two new parallel constrained optimization algorithms based on the variables distribution (PVD) have been presented. The first one consists of a parallel sequential quadratic programming approach for the case of block-separable constraints. This is the first PVD-type method whose convergence has been established for nonconvex feasible sets. The second proposed algorithm employs approximate projected gradient directions for the case of general (inseparable) convex constraints. The use of inexact directions is of particular relevance when the projection operation is computationally costly.

References

1. D. Bertsekas, *Nonlinear Programming*, Athena Scientific: Belmont, Massachusetts, 1995.
2. D. Bertsekas and J. Tsitsiklis, *Parallel and Distributed Computation*, Prentice-Hall, Inc.: Englewood Cliffs, New Jersey, 1989.
3. J. Bonnans, J. Gilbert, C. Lemaréchal, and C. Sagastizábal, *Optimisation Numérique: Aspects théoriques et pratiques*, Springer-Verlag: Berlin, 1997.
4. J. Conway and N. Sloane, *Sphere Packings, Lattices and Groups*, Springer-Verlag: New York, 1988.
5. R.W. Cottle, J.-S. Pang, and R.E. Stone, *The Linear Complementarity Problem*, Academic Press: New York, 1992.
6. M. Ferris and O. Mangasarian, “Parallel variable distribution,” *SIAM Journal on Optimization*, vol. 4, pp. 815–832, 1994.
7. M. Fukushima, “Parallel variable transformation in unconstrained optimization,” *SIAM Journal on Optimization*, vol. 8, pp. 658–672, 1998.
8. N. Gould, “On practical conditions for the existence and uniqueness of solutions to the general equality quadratic programming problem,” *Mathematical Programming*, vol. 32, pp. 90–99, 1985.
9. S.-P. Han, “A globally convergent method for nonlinear programming,” *Journal of Optimization Theory and Applications*, vol. 22, pp. 297–309, 1977.

10. S.-P. Han, "Optimization by updated conjugate subspaces," in *Numerical Analysis*, No. 140, Pitman Research Notes in Mathematics. Longman Scientific and Technical: Burnt Mill, England, 1986, pp. 82–97.
11. Y.Y. Lin and J.-S. Pang, "Iterative methods for large convex quadratic programs: A survey," *SIAM Journal on Control and Optimization*, vol. 25, pp. 383–411, 1987.
12. C.-S. Liu and C.-H. Tseng, "Parallel synchronous and asynchronous space-decomposition algorithms for large-scale minimization problems," *Computational Optimization and Applications*, vol. 17, pp. 85–107, 2000.
13. O. Mangasarian, *Nonlinear Programming*, McGraw-Hill: New York, 1969.
14. O. Mangasarian, "Parallel gradient distribution in unconstrained optimization," *SIAM Journal on Control and Optimization*, vol. 33, pp. 1916–1925, 1995.
15. O. Mangasarian and R. De Leone, "Error bounds for strongly convex programs and (super)linearly convergent iterative schemes for the least 2-norm solution of linear programs," *Applied Mathematics and Optimization*, vol. 17, pp. 1–14, 1988.
16. J.M. Martínez, "Two-phase model algorithm with global convergence for nonlinear programming," *Journal of Optimization Theory and Applications*, vol. 96, pp. 397–436, 1998.
17. D. Mayne and E. Polak, "A superlinearly convergent algorithm for constrained optimization problems," *Mathematical Programming Study*, vol. 16, pp. 45–61, 1982.
18. J.-S. Pang, "Error bounds in mathematical programming," *Mathematical Programming*, vol. 79, pp. 299–332, 1997.
19. M. Solodov, "New inexact parallel variable distribution algorithms," *Computational Optimization and Applications*, vol. 7, pp. 165–182, 1997.
20. M. Solodov, "On the convergence of constrained parallel variable distribution algorithms," *SIAM Journal on Optimization*, vol. 8, pp. 187–196, 1998.
21. P. Tseng, "Dual coordinate ascent methods for non-strictly convex minimization," *Mathematical Programming*, vol. 59, pp. 231–248, 1993.
22. E. Yamakawa and M. Fukushima, "Testing parallel variable transformation," *Computational Optimization and Applications*, vol. 13, pp. 253–274, 1999.