

Global Linear Convergence of an Augmented Lagrangian Algorithm to Solve Convex Quadratic Optimization Problems

Frédéric Delbos

*Institut Français du Pétrole,
1 & 4 Avenue de Bois-Préau, 92852 Rueil-Malmaison, France*

J. Charles Gilbert

*INRIA Rocquencourt, BP 105,
78153 Le Chesnay Cedex, France
jean-charles.gilbert@inria.fr*

This contribution is dedicated to Claude Lemaréchal, a friend and a mountaineering companion of the second author, on the occasion of his sixtieth birthday

Received October 31, 2003

We consider an augmented Lagrangian algorithm for minimizing a convex quadratic function subject to linear inequality constraints. Linear optimization is an important particular instance of this problem. We show that, provided the augmentation parameter is large enough, the constraint value converges *globally* linearly to zero. This property is viewed as a consequence of the proximal interpretation of the algorithm and of the global radial Lipschitz continuity of the reciprocal of the dual function subdifferential. This Lipschitz property is itself obtained by means of a lemma of general interest, which compares the distances from a point in the positive orthant to an affine space, on the one hand, and to the polyhedron given by the intersection of this affine space and the positive orthant, on the other hand. No strict complementarity assumption is needed. The result is illustrated by numerical experiments and algorithmic implications, including complexity issues, are discussed.

Keywords: Augmented Lagrangian, convex quadratic optimization, distance to a polyhedron, error bound, global linear convergence, iterative complexity, linear constraints, proximal algorithm

1991 Mathematics Subject Classification: 49M29, 65K05, 90C05, 90C06, 90C20

1. Introduction

Convex quadratic programs (QP) arise in their own right and as subproblems in some numerical algorithms to solve optimization problems. On the one hand, since no strictly convex assumption is made, the important family of linear optimization problems, with a zero quadratic term in their objective, enters this framework. On the other hand, the SQP algorithm decomposes a regularized constrained least-squares problem into a sequence of strictly convex QP's (see [10, 5] for an example in reflection tomography, which partly motivates this study; see [19, 2] for recent books describing the SQP algorithm). Finding efficient algorithms to solve this basic multi-faceted problem in all possible situations is an objective that has been pursued for decades (see for example the already 20 year old survey on quadratic programming in [20] and the monographs on interior point methods

in [15, 7, 18, 30, 14, 29, 31, 2]).

The convex quadratic problem we consider is written

$$\begin{cases} \min_x \frac{1}{2} x^\top Q x + q^\top x \\ l \leq Cx \leq u, \end{cases} \quad (1)$$

where Q is an $n \times n$ positive semi-definite symmetric matrix, $q \in \mathbb{R}^n$, C is $m \times n$, and the m -dimensional vectors l and u satisfy $l < u$ (i.e., $l_i < u_i$ for all indices $i = 1, \dots, m$) and may have infinite components. With lower and upper bounds, problem (1) is close to what is actually implemented in numerical codes (see [6] for an example). We have not included linear equality constraints, like $Ax = b$, in (1) to make the presentation simple, but such constraints can be expressed like in (1) by adding two inequalities $Ax \leq b$ and $-Ax \leq -b$. Therefore, our analysis covers problems with linear equality constraints. Note that, since Q may be zero, (1) also models linear optimization.

The method that we further explore in this paper fits into the class of dual approaches, since it is essentially the augmented Lagrangian (AL) method of Hestenes [11] and Powell [22] that is applied to the convex QP (1). This algorithm can be implemented in such a way that it does not require any matrix factorization. It is therefore appropriate when the problem is so large that such a factorization is impracticable or too much time consuming. This is a motivation for using the AL algorithm when the optimization problem deals with systems governed by partial differential equations [8]. In that case, however, a good preconditioner for the unavoidable conjugate gradient iterations must be available.

The version of the AL algorithm we analyze is defined on an equivalent form of (1) obtained by introducing an auxiliary variable $y \in \mathbb{R}^m$ [10]:

$$\begin{cases} \min_x \frac{1}{2} x^\top Q x + q^\top x \\ y = Cx \\ l \leq y \leq u. \end{cases} \quad (2)$$

The algorithm generates a sequence $\{\lambda_k\} \subset \mathbb{R}^m$ converging to some optimal multiplier associated with the equality constraint of (2). At each iteration, an auxiliary bound constrained QP has to be solved, so that the approach can be viewed as transforming (1) into a sequence of bound constrained convex quadratic subproblems. Two facts contribute to the possible success of this method. First, a bound constraint QP is much easier to solve than (1), which has general linear constraints (see [17, 9] and the references therein). Second, because of its dual and constraint convergence, the AL algorithm usually identifies the active constraints of (1) in a finite number of iterations. Since often these constraints are also the active constraints of the subproblems close to the solution, the combinatorial aspect of the bound constrained QP's rapidly decreases in intensity as the convergence progresses (and usually disappears after finitely many AL iterations). This reasoning is valid, for instance, when Q is positive definite and strict complementarity holds at the solution.

The AL algorithm also generates primal iterates $(x_k, y_k) \in \mathbb{R}^n \times \mathbb{R}^m$ and is controlled by the convergence of the constraint values to zero: if $\|y_k - Cx_k\|$ is less than a given tolerance, optimality can be considered to be reached. The algorithm is also driven by a so-called augmentation parameter r_k , whose role on the speed of convergence is major. This paper essentially shows that, provided r_k is larger than a certain positive threshold,

the convergence of the constraint norm to zero is *globally* linear, meaning that *at each iteration* the constraint norm decreases by a factor uniformly less than one. This property makes predictable the number of iterations to converge to a given precision and offers a possibility to study the global iterative complexity of the algorithm.

The paper is organized as follows. In Section 2, the AL algorithm under investigation is presented with the appropriate level of details. In Section 3, we give the tools from convex analysis that are useful for the study of the method. We already set out some of the properties of the algorithm. This section also gives a lemma of general interest, which compares the distances from a point in the positive orthant to an affine space, on the one hand, and to the polyhedron given by the intersection of this affine space and the positive orthant, on the other hand. Section 4 deals with the global linear convergence of the AL algorithm. It starts by showing a global error bound for the dual solution set, in terms of the subgradient of the dual function. The global linear convergence is then seen as an easy consequence of this property. We conclude in Section 5 by relating some numerical experiments on a seismic tomography problem and by a discussion on algorithmic consequences.

Notation

We denote the Euclidean norm by $\|\cdot\|$. The distance associated with this norm is denoted by “dist”, $B := \{x : \|x\| \leq 1\}$ is the closed unit ball, and $\partial B := \{x : \|x\| = 1\}$ is the unit sphere. We note $\mathbb{R} := \mathbb{R} \cup \{-\infty, +\infty\}$. The nonnegative orthant of \mathbb{R}^n is denoted by $\mathbb{R}_+^n := \{x \in \mathbb{R}^n : x \geq 0\}$, while $\mathbb{R}_{++}^n := \{x \in \mathbb{R}^n : x > 0\}$. The null space and range space of a matrix A are respectively denoted by $N(A)$ and $R(A)$. We write $A \succcurlyeq 0$ [resp. $A \succ 0$] to indicate that a symmetric matrix A is positive semi-definite [resp. positive definite].

Let E be a finite dimensional Euclidean space. The indicator function of a set $S \subset E$ is denoted by \mathcal{I}_S (this is the function that vanishes on S and takes the value $+\infty$ outside S). The domain of a function $f : E \rightarrow \mathbb{R} \cup \{+\infty\}$ is defined by $\text{dom } f := \{x \in E : f(x) < +\infty\}$ and its epigraph by $\text{epi } f := \{(x, \alpha) \in E \times \mathbb{R} : f(x) \leq \alpha\}$. As in [12], $\text{Conv}(E)$ is the set of functions $f : E \rightarrow \mathbb{R} \cup \{+\infty\}$ that are convex (epi f is convex), proper (epi $f \neq \emptyset$), and closed (epi f is closed). The subdifferential at $x \in E$ of a proper convex function $f : E \rightarrow \mathbb{R} \cup \{+\infty\}$ is denoted by $\partial f(x)$. We denote by $N_C(x)$ the normal cone at x to a convex set $C \subset E$. The orthogonal projection of a point x onto a nonempty closed convex set C is denoted by $P_C(x)$.

2. An AL algorithm to solve the QP

We **assume throughout that problem (2) has a solution** and denote by \mathcal{S}_P the set of its solutions (\bar{x}, \bar{y}) . The projections of \mathcal{S}_P onto \mathbb{R}^n and \mathbb{R}^m are respectively denoted by $\mathcal{S}_P^x := \{\bar{x} \in \mathbb{R}^n : (\bar{x}, \bar{y}) \in \mathcal{S}_P \text{ for some } \bar{y} \in \mathbb{R}^m\}$ (this is also $\{\bar{x} \in \mathbb{R}^n : (\bar{x}, C\bar{x}) \in \mathcal{S}_P\}$, the solution set of (1)) and $\mathcal{S}_P^y := \{\bar{y} \in \mathbb{R}^m : (\bar{x}, \bar{y}) \in \mathcal{S}_P \text{ for some } \bar{x} \in \mathbb{R}^n\}$. Since the constraints of (2) are qualified, there exist optimal multipliers, which certainly implies that the affine subspace

$$\Lambda := \{\lambda \in \mathbb{R}^m : C^\top \lambda \in q + R(Q)\} \tag{3}$$

is nonempty. Note that $\Lambda = \mathbb{R}^m$ if $Q \succ 0$.

The augmented Lagrangian is obtained by dualizing the equality constraint of (2). It is the function $\ell_r : (x, y, \lambda) \in \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R}^m \mapsto \mathbb{R}$, defined by

$$\ell_r(x, y, \lambda) = \frac{1}{2} x^\top Q x + q^\top x + \lambda^\top (y - Cx) + \frac{r}{2} \|y - Cx\|^2, \quad (4)$$

where $r \geq 0$ is called the *augmentation parameter* (see [1, 11, 22]).

We can now give a precise statement of the AL algorithm we study in this paper, which is basically the method of Hestenes [11] and Powell [22] applied to (2).

AL algorithm to solve (1)

Initialization: choose $\lambda_0 \in \mathbb{R}^m$ and $r_0 > 0$.

Repeat for $k = 0, 1, 2, \dots$

1. Solve:

$$\min_{\substack{x \\ l \leq y \leq u}} \ell_{r_k}(x, y, \lambda_k). \quad (5)$$

Denote a solution by (x_{k+1}, y_{k+1}) .

2. Update the multiplier

$$\lambda_{k+1} = \lambda_k + r_k (y_{k+1} - Cx_{k+1}). \quad (6)$$

3. Stop if $y_{k+1} \simeq Cx_{k+1}$.

4. Choose a new augmentation parameter: $r_{k+1} > 0$.

This algorithm deserves some comments.

1. Under the sole assumption that problem (1) has a solution, the QP in step (5) has also a solution. This fact is clarified in Proposition 3.3. This solution is not necessary unique however.
2. Even though (x_{k+1}, y_{k+1}) is not uniquely determined as a solution to (5), $y_{k+1} - Cx_{k+1}$ is independent of that solution, so that the multipliers λ_k are unambiguously generated.
3. The augmentation parameter r_k can change from iteration to iteration, but the same value must be used in the AL minimized in step 2 and in the multiplier update in step 2. If the “step-size” in (6) is different from r_k (with the aim at minimizing better the dual function, as in [21, Section 4.2] for example), several properties of the AL algorithm may no longer hold, such as the finite identification of active constraints (in the presence of strict complementarity) and the global linear convergence of Section 4.
4. The larger are the augmentation parameters r_k , the faster is the convergence. The only limitation on a large value for r_k comes from the ill-conditioning that such a value induces in the AL and the resulting difficulty in solving (5). Actually, it is clear from the structure of the AL in (4) that a large r gives priority to the restoration of the equality constraint, leaving aside the minimization of the Lagrangian (whose role is to provide optimality).

In comparison with an interior point method, which faces the combinatorial aspect of (1) by transforming the problem into a sequence of linear systems, the AL algorithm goes around this difficulty by transforming a general QP into a sequence of QP’s with simple bounds, which are easier to solve. Indeed, a number of efficient algorithms are available for

dealing with the bound constraints on the AL in step 2. A possibility would be to minimize first analytically ℓ_r in y and then to minimize the resulting function in x . Unfortunately this function of x , which is the AL associated with the inequality constrained QP (1) [24, 26], has a combinatorial structure (it contains maxima) that is not easier to deal with than the direct numerical minimization of ℓ_r in (x, y) with bounds on y . In our code QPAL [6], used for the numerical experiments of Section 5 and in [5], we have adapted to the (x, y) structure of problem (2) an active set strategy together with the gradient projection algorithm and conjugate gradient iterations on the activated faces (see [17, 9] and the references therein).

As opposed to standard (non shifted) interior point methods, whose elementary linear systems have an exploding condition number, the AL algorithm does not require the penalty parameter r_k to go to infinity. Actually, any sequence $\{r_k\}$ that remains bounded away from zero guarantees the convergence, even though large values speed it up [28]. Therefore, the bound constrained QP's in step 2 can be maintained reasonably well conditioned, keeping satisfactory the efficiency of a conjugate gradient based solver. This remark reinforces the viewpoint that considers the AL algorithm as a method suitable for large problems.

3. Convex analysis tools

3.1. Duality

As a dual function associated with problem (2), we use the one obtained by dualizing its equality constraints. It is the function $\delta : \mathbb{R}^m \rightarrow \mathbb{R} \cup \{+\infty\}$ defined by

$$\lambda \mapsto \delta(\lambda) := - \inf_{\substack{x \\ y \in [l, u]}} \left(\frac{1}{2} x^\top Q x + q^\top x + \lambda^\top (y - Cx) \right). \quad (7)$$

Clearly $\delta \in \overline{\text{Conv}}(\mathbb{R}^m)$ (it takes a finite value, for instance, when λ is an optimal multiplier associated with the equality constraint of (2)).

For a given $\lambda \in \mathbb{R}^m$, (x_λ, y_λ) is a solution to the *Lagrange problem*, the minimization problem in (7), if and only if $x_\lambda \in X_\lambda$ and $y_\lambda \in Y_\lambda$, where X_λ is the affine space

$$X_\lambda := \{x \in \mathbb{R}^n : Qx = C^\top \lambda - q\} \quad (8)$$

and Y_λ is the Cartesian product of the following intervals

$$(Y_\lambda)_i = \begin{cases} [u_i, u_i] & \text{if } \lambda_i < 0 \\ [l_i, u_i] & \text{if } \lambda_i = 0 \\ [l_i, l_i] & \text{if } \lambda_i > 0. \end{cases} \quad (9)$$

These intervals, with their possible infinite bounds, have to be understood in a broad sense: for example, $[l_i, u_i]$ is the interval $] -\infty, u_i]$ if $l_i = -\infty$ and u_i is finite, $[l_i, l_i]$ is the empty set if $l_i = -\infty$, etc. Since the Lagrange problem is always feasible and since a feasible convex quadratic problem has a solution if and only if it is bounded (see [2, Theorem 17.1] for example), the domain of δ is the set of λ 's for which the Lagrange problem has a solution. Therefore $\text{dom } \delta$ can be written as the nonempty polyhedron

$$\text{dom } \delta = \{\lambda \in \mathbb{R}^m : X_\lambda \neq \emptyset, Y_\lambda \neq \emptyset\} = \mathbb{R}_{l, u}^m \cap \Lambda,$$

where we used the fact that $X_\lambda \neq \emptyset$ if and only if $\lambda \in \Lambda$ and the notation

$$\mathbb{R}_{l,u}^m := \{\lambda \in \mathbb{R}^m : \lambda_i \leq 0 \text{ if } l_i = -\infty, \lambda_i \geq 0 \text{ if } u_i = +\infty\}.$$

Observe finally that the multivalued function $\lambda \mapsto -Y_\lambda$ is monotone: for λ and $\lambda' \in \mathbb{R}^m$, and for $y_\lambda \in Y_\lambda$ and $y_{\lambda'} \in Y_{\lambda'}$, there holds

$$-(y_{\lambda'} - y_\lambda)^\top (\lambda' - \lambda) \geq 0. \quad (10)$$

Let Q^\dagger be the pseudo-inverse of Q and take the notation

$$H := CQ^\dagger C^\top \quad \text{and} \quad v := CQ^\dagger q. \quad (11)$$

Let $\lambda \in \text{dom } \delta$. Then, the minimization in x of the Lagrangian in (7) is achieved by any $x_\lambda \in X_\lambda$, hence satisfying $Qx_\lambda + q = C^\top \lambda$. This minimizer is not necessarily unique and we can take the minimum norm minimizer $x_\lambda^\dagger := Q^\dagger(C^\top \lambda - q)$. After substitution in the function minimized in (7) and the use of $Q^\dagger Q Q^\dagger = Q^\dagger$, one gets

$$\delta(\lambda) = \sup_{y \in [l,u]} \left(\frac{1}{2} \lambda^\top H \lambda - (v + y)^\top \lambda + \frac{1}{2} q^\top Q^\dagger q \right), \quad \text{for } \lambda \in \text{dom } \delta. \quad (12)$$

On its domain, the dual function δ is therefore the maximum of a finite number of convex quadratic functions (those defined by the argument of the supremum in (12) with the components of y set to l_i or u_i ; only the finite values of these bounds must be considered), which only differ by their slope at the origin (in particular, they have the same Hessian H).

Lemma 3.1. *The subdifferential of the dual function (7) is given at $\lambda \in \text{dom } \delta$ by*

$$\partial \delta(\lambda) = H\lambda - v - Y_\lambda + C(N(Q)).$$

Proof. We write δ as the sum of three convex functions. Let \tilde{l} and $\tilde{u} \in \mathbb{R}^m$ be chosen such that $\tilde{l} < \tilde{u}$, $\tilde{l}_i = l_i$ if l_i is finite, and $\tilde{u}_i = u_i$ if u_i is finite. Define \tilde{Y}_λ by formula (9) with l and u respectively replaced by \tilde{l} and \tilde{u} . Then the following *finite* value function $\tilde{\delta} \in \text{Conv}(\mathbb{R}^m)$ is identical to the right hand side of (12) on $\mathbb{R}_{l,u}^m$:

$$\tilde{\delta}(\lambda) = \sup_{\substack{y=(y_i)_{i=1}^m \\ y_i = \tilde{l}_i \text{ or } \tilde{u}_i}} \left(\frac{1}{2} \lambda^\top H \lambda - (v + y)^\top \lambda + \frac{1}{2} q^\top Q^\dagger q \right).$$

Clearly $\delta = \tilde{\delta} + \mathcal{I}_{\mathbb{R}_{l,u}^m} + \mathcal{I}_\Lambda$, so that Theorem 23.8 in [23] implies that for $\lambda \in \text{dom } \delta$:

$$\partial \delta(\lambda) = \partial \tilde{\delta}(\lambda) + \partial \mathcal{I}_{\mathbb{R}_{l,u}^m}(\lambda) + \partial \mathcal{I}_\Lambda(\lambda).$$

Equality holds above because $\mathcal{I}_{\mathbb{R}_{l,u}^m}$ and \mathcal{I}_Λ are polyhedral and because $\text{ri dom } \tilde{\delta} (= \mathbb{R}^m, \text{ri denotes the relative interior})$, $\text{dom } \mathcal{I}_{\mathbb{R}_{l,u}^m} = \mathbb{R}_{l,u}^m$, and $\text{dom } \mathcal{I}_\Lambda = \Lambda$ have a point in common (one in $\text{dom } \delta \neq \emptyset$).

To compute $\partial \tilde{\delta}(\lambda)$, we use Corollary VI.4.3.2 in [12] (conv denotes the convex hull):

$$\partial \tilde{\delta}(\lambda) = \text{conv} \left\{ H\lambda - v - y : y \in \tilde{Y}_\lambda \text{ and } (y_i = \tilde{l}_i \text{ or } \tilde{u}_i \text{ if } \lambda_i = 0) \right\} = H\lambda - v - \tilde{Y}_\lambda.$$

On the other hand, $\partial\mathcal{I}_{\mathbb{R}^m}(\lambda) = N_{\mathbb{R}^m}(\lambda)$, which is the set of vectors $\nu \in \mathbb{R}^m$ satisfying

$$\begin{cases} \nu_i \geq 0 & \text{when } \lambda_i = 0, l_i = -\infty, \text{ and } u_i \text{ is finite} \\ \nu_i \leq 0 & \text{when } \lambda_i = 0, l_i \text{ is finite, and } u_i = +\infty \\ \nu_i \in \mathbb{R} & \text{when } \lambda_i = 0, l_i = -\infty, \text{ and } u_i = +\infty \\ \nu_i = 0 & \text{when } \lambda_i \neq 0. \end{cases}$$

We deduce from this computation that for $\lambda \in \text{dom } \delta$

$$\tilde{Y}_\lambda - \partial\mathcal{I}_{\mathbb{R}^m}(\lambda) = Y_\lambda.$$

Finally $\partial\mathcal{I}_\Lambda(\lambda) = N_\Lambda(\lambda) = \{\mu \in \mathbb{R}^m : C^\top \mu \in R(Q)\}^\perp = C(N(Q))$. Adding the last three subdifferentials provides the formula of $\partial\delta(\lambda)$ given in the statement of the lemma. \square

Let us denote by \mathcal{S}_D the set of dual solutions:

$$\mathcal{S}_D := \{\bar{\lambda} \in \mathbb{R}^m : 0 \in \partial\delta(\bar{\lambda})\}.$$

Not surprisingly, this is a convex polyhedron, which can be described in the standard form. It will be useful to make this form explicit and we do so in Lemma 3.2 below. For this, we take a partition of $\{1, \dots, m\}$ into the index sets

$$\begin{aligned} I_l &:= \{i : \bar{y}_i = l_i \text{ for all } (\bar{x}, \bar{y}) \in \mathcal{S}_P\}, \\ J &:= \{i : l_i < \bar{y}_i < u_i \text{ for some } (\bar{x}, \bar{y}) \in \mathcal{S}_P\}, \\ I_u &:= \{i : \bar{y}_i = u_i \text{ for all } (\bar{x}, \bar{y}) \in \mathcal{S}_P\}. \end{aligned} \tag{13}$$

We also introduce the orthant face \mathcal{O} and the affine subspace \mathcal{A}

$$\begin{aligned} \mathcal{O} &:= \{\lambda \in \mathbb{R}^m : \lambda_{I_l} \geq 0, \lambda_J = 0, \lambda_{I_u} \leq 0\}, \\ \mathcal{A} &:= \{\lambda \in \mathbb{R}^m : C^\top \lambda = Q\bar{x} + q\}. \end{aligned} \tag{14}$$

In the definition of \mathcal{A} , \bar{x} is an arbitrary primal solution. We have not made this dependence explicit in the symbol of the set since, as shown in the proof of the next lemma, \mathcal{A} does not depend on the choice of $\bar{x} \in \mathcal{S}_P^x$.

Lemma 3.2. *The set of dual solutions \mathcal{S}_D can be written as the intersection*

$$\mathcal{S}_D = \mathcal{O} \cap \mathcal{A}.$$

Furthermore, for any $\bar{\lambda} \in \mathcal{S}_D$ and any $\bar{y} \in \mathcal{S}_P^y$, we have $\bar{y} \in Y_{\bar{\lambda}}$ and $H\bar{\lambda} = v + \bar{y} + C\bar{u}$ for some $\bar{u} \in N(Q)$.

Proof. Let $\ell(x, y, \lambda) = \frac{1}{2}x^\top Qx + q^\top x + \mathcal{I}_{[l, u]}(y) + \lambda^\top(y - Cx)$ be the Lagrangian function of the problem $\min_{(x, y)} \{\frac{1}{2}x^\top Qx + q^\top x + \mathcal{I}_{[l, u]}(y) : y = Cx\}$, which has the same dual function as problem (2). Since the constraint of this problem is qualified, $\bar{\lambda} \in \mathcal{S}_D$ if and only if $0 \in \partial_{(x, y)}\ell(\bar{x}, \bar{y}, \bar{\lambda})$, where (\bar{x}, \bar{y}) is an arbitrary primal solution. This can also be written $Q\bar{x} + q = C^\top \bar{\lambda}$ and $0 \in N_{[l, u]}(\bar{y}) + \bar{\lambda}$, which is equivalent to $\bar{\lambda} \in \mathcal{A}_{\bar{x}} \cap \mathcal{O}_{\bar{y}}$, where

$$\begin{aligned} \mathcal{A}_{\bar{x}} &:= \{\lambda \in \mathbb{R}^m : C^\top \lambda = Q\bar{x} + q\}, \\ \mathcal{O}_{\bar{y}} &:= \{\lambda \in \mathbb{R}^m : \lambda_i \geq 0 \text{ if } \bar{y}_i = l_i, \lambda_i = 0 \text{ if } l_i < \bar{y}_i < u_i, \lambda_i \leq 0 \text{ if } \bar{y}_i = u_i\}. \end{aligned}$$

By varying $(\bar{x}, \bar{y}) \in \mathcal{S}_P$, we see that $\mathcal{A}_{\bar{x}} = \mathcal{A}$ is independent of the chosen primal solution $\bar{x} \in \mathcal{S}_P^x$ and that $\bar{\lambda} \in \cap\{\mathcal{O}_{\bar{y}} : \bar{y} \in \mathcal{S}_P^y\} = \mathcal{O}$.

For proving the second part of the lemma, take $\bar{\lambda} \in \mathcal{S}_D$ and $(\bar{x}, \bar{y}) \in \mathcal{S}_P$. We have shown that $\bar{\lambda} \in \mathcal{A}_{\bar{x}} \cap \mathcal{O}_{\bar{y}}$. Actually, $\bar{\lambda} \in \mathcal{O}_{\bar{y}}$ is equivalent to $\bar{y} \in Y_{\bar{\lambda}}$. By $\bar{\lambda} \in \mathcal{A}_{\bar{x}}$, we have that $C^\top \bar{\lambda} = Q\bar{x} + q$. Multiplying to the left both sides of this equation by CQ^\dagger provides $H\bar{\lambda} = v + \bar{y} + C\bar{u}$, where $\bar{u} := (Q^\dagger Q - I)\bar{x} \in N(Q)$. \square

The fact observed in the proof above that the gradient of the criterion of the primal problem at a solution, here $Q\bar{x} + q$, is independent of the chosen solution is a property of general convex problems; see [16, 3]. This fact can also be deduced from the property that the subdifferential of a convex function (here the criterion of problem (1)) is constant on the relative interior of a set on which this function is constant (here the solution set).

3.2. Proximalilty

We will use the fundamental result of Rockafellar [25], according to which the AL algorithm of Section 2 is the proximal algorithm on the dual function δ . More precisely, the multiplier λ_{k+1} computed in step 2 of the AL algorithm is also the unique solution to

$$\inf_{\lambda \in \mathbb{R}^m} \left(\delta(\lambda) + \frac{1}{2r_k} \|\lambda - \lambda_k\|^2 \right). \quad (15)$$

The same parameter $r_k > 0$ is used above and in (5). In addition, the optimal value of this problem is the opposite of the optimal value of problem (5). The optimality conditions of problem (15) can be written $0 \in \partial\delta(\lambda_{k+1}) + (\lambda_{k+1} - \lambda_k)/r_k$. Using (6), we see that:

$$Cx_k - y_k \in \partial\delta(\lambda_k), \quad \forall k \geq 1. \quad (16)$$

Note that, since λ_{k+1} is uniquely determined as the solution to (15), this is also the case for $y_{k+1} - Cx_{k+1}$, even though x_{k+1} and y_{k+1} are not uniquely determined.

Let us now clarify the conditions ensuring that the augmented Lagrange problem (5) has a solution.

Proposition 3.3. *The following three properties are equivalent:*

- (i) $\text{dom } \delta \neq \emptyset$,
- (ii) *problem (1), with some (or any) finite shift of its finite bounds to make it feasible, has a solution,*
- (iii) *for some (or any) $r_k > 0$ and $\lambda_k \in \mathbb{R}^m$, problem (5) has a solution.*

(i) \Rightarrow (iii). Fix $r_k > 0$ and $\lambda_k \in \mathbb{R}^m$ (not necessarily the k th iterate). Since $\text{dom } \delta \neq \emptyset$, the optimal value of (15) is finite, so that the optimal value of problem (5) is also finite. As a feasible bounded convex quadratic problem, (5) must have a solution [2, Theorem 17.1].

[(iii) \Rightarrow (ii)] We proceed by contradiction. Suppose that \tilde{l} and $\tilde{u} \in \bar{\mathbb{R}}^m$ are such that $\tilde{l} < \tilde{u}$, $\tilde{l}_i = -\infty$ iff $l_i = -\infty$, $\tilde{u}_i = +\infty$ iff $u_i = +\infty$, and $[\tilde{l}, \tilde{u}] \cap R(C) \neq \emptyset$ (these bounds \tilde{l} and \tilde{u} result from a finite shift of the finite bounds of (1) that makes this problem feasible) and assume that the feasible problem $\min\{f(x) : \tilde{l} \leq Cx \leq \tilde{u}\}$, where $f(x) := (1/2)x^\top Qx + q^\top x$, has no solution. Then, there exists a sequence $\{x_j\}$ such

that $Cx_j \in [\tilde{l}, \tilde{u}]$ and $f(x_j) \rightarrow -\infty$ when $j \rightarrow \infty$ (this is because a bounded feasible convex quadratic problem has a solution). Let $y_j := P_{[l, u]}(Cx_j)$ be the projection of Cx_j onto $[l, u]$. Then $\|y_j - Cx_j\| \leq m^{1/2}\|y_j - Cx_j\|_\infty \leq m^{1/2}\gamma$, where $\gamma := \max(\|\tilde{l} - l\|_\infty, \|\tilde{u} - u\|_\infty)$ (these norms are taken on the finite components of l and u), and (x_j, y_j) is feasible for problem (5). On the other hand, for an arbitrary $r_k > 0$ and $\lambda_k \in \mathbb{R}^m$, $\ell_{r_k}(x_j, y_j, \lambda_k) \leq f(x_j) + m^{1/2}\gamma\|\lambda_k\| + (r_k/2)m\gamma^2 \rightarrow -\infty$ when $j \rightarrow \infty$. Therefore problem (5) has no solution.

[(ii) \Rightarrow (i)] Let \tilde{l} and \tilde{u} be some finite shifts of the finite bounds of problem (1), such that the problem $\min\{f(x) : \tilde{l} \leq Cx \leq \tilde{u}\}$, with f as in the previous paragraph, has a solution, \tilde{x} say. Since its constraints are qualified, there exist $\tilde{\lambda}^l$ and $\tilde{\lambda}^u$ such that $Q\tilde{x} + q = C^\top(\tilde{\lambda}^l - \tilde{\lambda}^u)$, $\tilde{\lambda}^l \geq 0$, $\tilde{\lambda}^u \geq 0$, $\tilde{\lambda}_i^l = 0$ if $l_i = -\infty$, and $\tilde{\lambda}_i^u = 0$ if $u_i = +\infty$. It is easy to check that $\tilde{\lambda}^l - \tilde{\lambda}^u \in \mathbb{R}_{l, u}^m \cap \Lambda = \text{dom } \delta$. \square

If the original quadratic problem (1) has a solution, condition (ii) above holds (without having to shift the bounds), so that the augmented Lagrange problem (5) has a solution.

3.3. Projection onto a convex polyhedron

This section gives two lemmas related to the projection onto a convex polyhedron. The first lemma has a general interest. It compares the distance from a point x in the positive orthant to a convex polyhedron \mathcal{X} defined in the standard form and the distance from x to the underlying affine space \mathcal{A} . It is claimed that the second distance is bounded below by a positive constant (independent of $x \geq 0$) times the first one. Of course, since $\mathcal{X} \subset \mathcal{A}$, $\text{dist}(x, \mathcal{A}) \leq \text{dist}(x, \mathcal{X})$.

Lemma 3.4. *Let A be an $m \times n$ matrix and $b \in \mathbb{R}^m$. Consider the affine subspace \mathcal{A} and the convex polyhedron \mathcal{X} defined by*

$$\mathcal{A} := \{x \in \mathbb{R}^n : Ax = b\} \quad \text{and} \quad \mathcal{X} := \{x \in \mathbb{R}^n : Ax = b, x \geq 0\}.$$

These sets are supposed to be nonempty. Then, there exists a constant $\gamma > 0$ such that

$$\forall x \in \mathbb{R}_+^n, \quad \text{dist}(x, \mathcal{A}) \geq \gamma \text{dist}(x, \mathcal{X}).$$

Proof. *First stage: reformulation of the statement of the lemma.* By using the triangle inequality, it is easy to see that the conclusion of the lemma is equivalent to claiming that

$$\exists \gamma' > 0, \quad \forall x \in \mathbb{R}_+^n, \quad \|x - P_{\mathcal{A}}(x)\| \geq \gamma' \|P_{\mathcal{A}}(x) - P_{\mathcal{X}}(x)\|, \quad (17)$$

This inequality suggests that a certain function (whose value is the left hand side of the inequality in (17) divided by the factor of γ' in the right hand side) has a positive slope. This is the strategy we follow to establish (17).

Instead of testing the validity of the inequality in (17) for any $x \in \mathbb{R}_+^n$, we consider all the possible points x^0 that are projections onto \mathcal{X} of points in the positive orthant and reformulate (17). For a point $x^0 \in \mathcal{X}$ and a direction $d \in N_{\mathcal{X}}(x^0) \cap N(\mathcal{A}) \cap \partial B$, let us introduce the function

$$\varphi : \alpha \in \mathbb{R}_{++} \mapsto \varphi(\alpha) := \inf \{ \|A^\top y\| : y \in \mathbb{R}^m, x^0 + \alpha d + A^\top y \geq 0 \}. \quad (18)$$

This function depends on the choice of x^0 and d , but we do not mention this dependence to keep the notation light. Let us show that (17) holds if

$$\exists \gamma' > 0, \quad \forall x^0 \in \mathcal{X}, \quad \forall d \in N_{\mathcal{X}}(x^0) \cap N(A) \cap \partial B, \quad \forall \alpha > 0, \quad \varphi(\alpha) \geq \gamma' \alpha. \quad (19)$$

Let $x \in \mathbb{R}_+^n$, $x^1 := P_{\mathcal{A}}(x)$, and $x^0 := P_{\mathcal{X}}(x)$. One can assume that $x^1 \neq x^0$ (since otherwise the inequality in (17) is trivially satisfied). Set $d := (x^1 - x^0)/\|x^1 - x^0\|$. It is clear that $d \in N(A) \cap \partial B$ (note that both x^1 and $x^0 \in \mathcal{A}$). To show that $d \in N_{\mathcal{X}}(x^0)$, observe that $x^1 := P_{\mathcal{A}}(x)$ implies that $x = x^1 + A^\top y = x^0 + \alpha d + A^\top y$ for $\alpha := \|x^1 - x^0\| > 0$ and a certain $y \in \mathbb{R}^m$. Then, for all $z \in \mathcal{X}$, there holds

$$d^\top(z - x^0) = \frac{1}{\alpha}(x - x^0 - A^\top y)^\top(z - x^0) = \frac{1}{\alpha}(x - x^0)^\top(z - x^0) \leq 0.$$

We have used the fact that $z - x^0 \in N(A)$ and that $x^0 := P_{\mathcal{X}}(x)$ to get the last inequality. This shows that $d \in N_{\mathcal{X}}(x^0)$. Using (19) and the fact that $x = x^0 + \alpha d + A^\top y \geq 0$, we see that the inequality in (17) holds:

$$\|x - x^1\| = \|A^\top y\| \geq \varphi(\alpha) \geq \gamma' \alpha = \gamma' \|x^1 - x^0\|.$$

The claim (19) can be simplified. Observe that $\varphi(\alpha) \geq 0$, that $\varphi(0) = 0$, and that $\varphi(t\alpha) \leq t\varphi(\alpha)$ when $\alpha \geq 0$ and $t \in]0, 1]$. To prove this last property of φ , assume that $\varphi(\alpha) < \infty$ (otherwise, there is nothing to show). Then take $\varepsilon > 0$ and $y \in \mathbb{R}^m$ such that $x^0 + \alpha d + A^\top y \geq 0$ and $\|A^\top y\| \leq \varphi(\alpha) + \varepsilon$. Since $x^0 \geq 0$, there holds $0 \leq (1-t)x^0 + t(x^0 + \alpha d + A^\top y) = x^0 + t\alpha d + A^\top(ty)$ and therefore $\varphi(t\alpha) \leq \|A^\top(ty)\| \leq t\varphi(\alpha) + \varepsilon$. Since $\varepsilon > 0$ is arbitrary, there holds $\varphi(t\alpha) \leq t\varphi(\alpha)$. Now, this property of φ implies that $\alpha \in \mathbb{R}_{++} \mapsto \varphi(\alpha)/\alpha$ is nondecreasing. Therefore, we have reduced the problem to showing that

$$\exists \gamma' > 0, \quad \forall x^0 \in \mathcal{X}, \quad \forall d \in N_{\mathcal{X}}(x^0) \cap N(A) \cap \partial B, \quad \varphi'(0; 1) \geq \gamma', \quad (20)$$

where $\varphi'(0; 1)$ denotes the right derivative of φ at zero.

Second stage: control of the decomposition of the normal directions. Consider a point $x^0 \in \mathcal{X}$ having a unitary normal direction in the null space of A , say $d \in N_{x^0}(\mathcal{X}) \cap N(A) \cap \partial B$. Define $I := I(x^0) := \{i : x_i^0 = 0\}$ and $J := J(x^0) := \{i : x_i^0 > 0\}$. These directions d are characterized by the conditions

$$d = A^\top z - r, \quad r_I \geq 0, \quad r_J = 0, \quad Ad = 0, \quad \text{and} \quad \|d\| = 1, \quad (21)$$

for some vectors $z \in \mathbb{R}^m$ and $r \in \mathbb{R}^n$. The decomposition of d in $A^\top z - r$ as above is not necessarily unique. It will be useful to identify a decomposition that provides the smallest value to $\|A^\top z\|$, which is therefore a solution to

$$\begin{cases} \min_{(z,r)} \frac{1}{2} \|A^\top z\|^2 \\ A^\top z - r = d \\ r_I \geq 0 \\ r_J = 0. \end{cases}$$

It is easy to show that this problem has a solution, which is characterized by (21) and

$$A(A^\top z - s) = 0, \quad s_I \geq 0, \quad \text{and} \quad s_I^\top r_I = 0, \quad (22)$$

for some vector $s \in \mathbb{R}^n$.

Let us show that

$$\max_{x^0 \in \mathcal{X}} \sup_{\substack{d \in N_{\mathcal{X}}(x^0) \\ Ad=0 \\ \|d\|=1}} \min_{\substack{(z,r) \in \mathbb{R}^m \times \mathbb{R}^n \\ A^\top z - r = d \\ r_{I(x^0)} \geq 0 \\ r_{J(x^0)} = 0}} \|A^\top z\| < +\infty. \quad (23)$$

We see on (21) that two points $x^0 \in \mathcal{X}$ having the same index set I have the same normal cone. Therefore, the point $x^0 \in \mathcal{X}$ intervenes in (23) only through its index sets I and J . Since there is a finite number of such sets, one can fix x^0 , hence I and J . Let us continue by contradiction, assuming that there exists a sequence $\{(d^k, z^k, r^k, s^k)\}$ such that $Ad^k = 0$, $\|d^k\| = 1$, $A^\top z^k - r^k = d^k$, $r_I^k \geq 0$, $r_J^k = 0$, $A(A^\top z^k - s^k) = 0$, $s_I^k \geq 0$, $(s_I^k)^\top r_I^k = 0$, and $\|A^\top z^k\| \rightarrow \infty$. Extracting a subsequence if necessary, it can be assumed that $A^\top z^k / \|A^\top z^k\| \rightarrow A^\top \bar{z}$, a vector of unit norm. Since $\{d^k\}$ is bounded, the identity $A^\top z^k - r^k = d^k$ shows that $r^k / \|A^\top z^k\|$ converges to $\bar{r} := A^\top \bar{z}$. Multiplying the identity $A(A^\top z^k - s^k) = 0$ by \bar{z} , one finds for sufficiently large k

$$0 = \bar{z}^\top AA^\top z^k - \bar{r}^\top s^k = \bar{z}^\top AA^\top z^k,$$

because, when $\bar{r}_i > 0$, then $i \in I$ and, for all sufficiently large k , $r_i^k > 0$, so that $s_i^k = 0$. Dividing the right hand side by $\|A^\top z^k\|$ and taking the limit, one would find $A^\top \bar{z} = 0$, which provides the expected contradiction.

Third stage: lower bound for $\varphi'(0; 1)$ and conclusion. Let us introduce $v : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$, the value function of the problem

$$\begin{cases} \inf_y \|A^\top y\| \\ x^0 + A^\top y \geq 0, \end{cases} \quad (24)$$

which is the proper convex function defined by $v(p) := \inf\{\|A^\top y\| : x^0 + A^\top y \geq p\}$. Then, for fixed $x^0 \in \mathcal{X}$ and $d \in N_{\mathcal{X}}(x^0) \cap N(A) \cap \partial B$, $\varphi(\alpha)$ defined by (18) can be written $\varphi(\alpha) = v(-\alpha d)$. Therefore

$$\varphi'(0; 1) = v'(0; -d) \geq -g^\top d, \quad \forall g \in \partial v(0). \quad (25)$$

As for the subdifferential $\partial v(0)$, it is formed of the optimal multipliers associated with the constraint of (24), which are the g -parts of the pairs (g, u) satisfying

$$g \in u + N(A), \quad \|u\| \leq 1, \quad g \geq 0, \quad \text{and} \quad (x^0)^\top g = 0. \quad (26)$$

Let $d = A^\top z - r$ be a decomposition of d satisfying (21)-(22). If $A^\top z = 0$, then $g := -\alpha d = \alpha r$ is a subgradient of v at zero for any $\alpha \geq 0$ (the conditions (26) are satisfied with $u = 0$; recall that $d \in N(A)$), so that (25) shows that $\varphi'(0; 1) = +\infty$. If $A^\top z \neq 0$, then $g := (A^\top z - d) / \|A^\top z\| = r / \|A^\top z\|$ is a subgradient of v at zero (the conditions (26) are satisfied with $u = A^\top z / \|A^\top z\|$). Then (25) shows that $\varphi'(0; 1) \geq 1 / \|A^\top z\|$. Since for these decompositions, stage 2 of the proof has shown that $A^\top z$ is bounded, (20) holds and, consequently, the result is proven. \square

As shown by the following example, Lemma 3.4 no longer holds with all its generality when \mathcal{X} is the intersection of an affine space \mathcal{A} and an arbitrary closed convex cone.

Example 3.5. Let us introduce the following closed convex cone $K := \{x \in \mathbb{R}^3 : x_2 x_3 \geq x_1^2, x_2 \geq 0, x_3 \geq 0\}$, the 1×3 matrix $A := (0 \ 1 \ 0)$, and $b = 0 \in \mathbb{R}$. Define the affine space \mathcal{A} and its intersection with K by

$$\begin{aligned}\mathcal{A} &:= \{x \in \mathbb{R}^3 : Ax = b\} = \{x \in \mathbb{R}^3 : x_2 = 0\}, \\ \mathcal{X} &:= K \cap \mathcal{A} = \{x \in \mathbb{R}^3 : x_1 = x_2 = 0, x_3 \geq 0\}.\end{aligned}$$

Then the conclusion of Lemma 3.4 does not hold for these sets \mathcal{A} and \mathcal{X} . To see this, consider the points $x^t := (t, t^2, 1)$ for $t \downarrow 0$. Clearly $x^t \in K$, $P_{\mathcal{A}}(x^t) = (t, 0, 1)$, and $P_{\mathcal{X}}(x^t) = (0, 0, 1)$. Therefore $\|x^t - P_{\mathcal{A}}(x^t)\|/\|P_{\mathcal{A}}(x^t) - P_{\mathcal{X}}(x^t)\| = t$, which is not bounded away from zero. \square

Actually, it will be useful below to have the following relaxed version of Lemma 3.4. This one allows the projected point x not to belong to \mathbb{R}_+^n . This point must however be sufficiently close to the positive orthant with respect to its distance to \mathcal{X} .

Corollary 3.6. *Assume the framework defined in the statement of Lemma 3.4. Then, there exist two constants $\tau > 0$ and $\gamma > 0$ such that for all $x \in \mathbb{R}^n$,*

$$\text{dist}(x, \mathbb{R}_+^n) \leq \tau \text{dist}(x, \mathcal{X}) \quad \implies \quad \text{dist}(x, \mathcal{A}) \geq \gamma \text{dist}(x, \mathcal{X}).$$

Proof. Let γ be the constant given by Lemma 3.4 and set

$$\tau := \frac{\gamma}{4(1 + \gamma)}.$$

Let x be such that $\text{dist}(x, \mathbb{R}_+^n) \leq \tau \text{dist}(x, \mathcal{X})$. To simplify the notation, let us define

$$x^0 := P_{\mathcal{X}}(x), \quad x^1 := P_{\mathcal{A}}(x), \quad \text{and} \quad \bar{x} := P_{\mathbb{R}_+^n}(x),$$

Using several times the triangle inequality, Lemma 3.4, the non-expansiveness of the projectors $P_{\mathcal{A}}$ and $P_{\mathcal{X}}$, and the definition of τ , one can write

$$\begin{aligned}\|x - x^1\| &\geq \|\bar{x} - P_{\mathcal{A}}(\bar{x})\| - \|P_{\mathcal{A}}(\bar{x}) - P_{\mathcal{A}}(x)\| - \|x - \bar{x}\| \\ &\geq \gamma \|\bar{x} - P_{\mathcal{X}}(\bar{x})\| - 2\|x - \bar{x}\| \\ &\geq \gamma \|\bar{x} - x^0\| - (2 + \gamma)\|x - \bar{x}\| \\ &\geq \gamma \|x - x^0\| - 2(1 + \gamma)\|x - \bar{x}\| \\ &\geq \gamma \|x - x^0\| - 2\tau(1 + \gamma)\|x - x^0\| \\ &= \frac{\gamma}{2} \|x - x^0\|.\end{aligned}$$

This is the expected inequality. \square

The following lemma will be also useful. If $I \subset \{1, \dots, n\}$, we denote by I^c the complementary set of I in $\{1, \dots, n\}$.

Lemma 3.7. *Let A be an $m \times n$ matrix, $b \in \mathbb{R}^m$, $I \subset \{1, \dots, n\}$, and $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex differentiable function. If \bar{x} is a solution to the problem*

$$\min \{\varphi(x) : Ax = b, x_I \geq 0, x_{I^c} = 0\},$$

then there is a subset of indices $J \subset \{1, \dots, n\}$, containing I , such that \bar{x} is also a solution to the problem

$$\min \{\varphi(x) : Ax = b, x_J \geq 0, x_{J^c} \leq 0\}.$$

Proof. The constraints of the first problem are affine, hence qualified. Therefore, there exist vectors $y \in \mathbb{R}^m$ and $s \in \mathbb{R}^n$ such that

$$\nabla\varphi(\bar{x}) + A^\top y + s = 0, \quad \bar{x}_I \geq 0, \quad s_I \leq 0, \quad s_I^\top \bar{x}_I = 0, \quad \bar{x}_{I^c} = 0.$$

Define

$$J := I \cup \{i \in I^c : s_i \leq 0\}.$$

Then

$$\begin{aligned} \nabla\varphi(\bar{x}) + A^\top y + s = 0, \quad \bar{x}_J \geq 0, \quad s_J \leq 0, \quad s_J^\top \bar{x}_J = 0, \\ \bar{x}_{J^c} \leq 0, \quad s_{J^c} \geq 0, \quad s_{J^c}^\top \bar{x}_{J^c} = 0. \end{aligned}$$

By convexity, these conditions suffice to show that \bar{x} is also a solution to the second problem. \square

4. Global linear convergence of the algorithm

The global linear convergence of the AL algorithm will be shown in Section 4.2 to be a consequence of the radial Lipschitz continuity of the multifunction $\partial\delta^{-1}$, the reciprocal of the subdifferential of the dual function (this argument is taken from [28]). The latter property is the subject of Section 4.1.

4.1. Global dual error bound

Two normed spaces E and F being given, a multifunction $T : E \rightarrow F$ is said to be *radially Lipschitz continuous at $x_0 \in E$ with constant $L \geq 0$* if for all $x \in E$ and all $y \in T(x)$, there holds $\text{dist}(y, T(x_0)) \leq L\|x - x_0\|$ (“dist” denotes here the distance associated with the norm of F). Consider the multifunction

$$\partial\delta^{-1} : g \in \mathbb{R}^m \mapsto \{\lambda \in \mathbb{R}^m : g \in \partial\delta(\lambda)\} \subset \mathbb{R}^m,$$

where δ is the dual function defined in (7). Clearly $\partial\delta^{-1}(0) = \mathcal{S}_D$, the set of dual solutions. Then $\partial\delta^{-1}$ is radially Lipschitz continuous at 0 with constant $L \geq 0$, if

$$\forall \lambda \in \mathbb{R}^m, \forall g \in \partial\delta(\lambda) : \text{dist}(\lambda, \mathcal{S}_D) \leq L\|g\|. \tag{27}$$

Such a property is sometimes called a *global error bound* for the dual solution set \mathcal{S}_D in terms of the dual function subgradient (see the review paper by Pang [21] and the contribution of Izmailov and Solodov [13]). In this section, we show that this property holds in a weaker form: λ has to stay at a bounded distance from \mathcal{S}_D (the Lipschitz constant L depends on this distance). Nevertheless, this property still has a global nature, since λ is not required to be close to \mathcal{S}_D and g is not required to be close to 0.

To show that this property is natural, consider first a quadratic problem with only equality constraints:

$$\begin{cases} \inf_x \frac{1}{2} x^\top Qx + q^\top x \\ Cx = b. \end{cases} \tag{28}$$

It is assumed that this problem is convex ($Q \succcurlyeq 0$) and has a solution. It is therefore feasible: $b \in R(C)$. Since the constraint is qualified, there exist optimal multipliers, which implies that the affine subspace Λ defined in (3) is nonempty.

Using the pseudo-inverse Q^\dagger of Q , the symmetric matrix $H \succcurlyeq 0$, and the vector v defined in (11), the dual function δ associated with problem (28) can be written

$$\delta(\lambda) = \begin{cases} \frac{1}{2} \lambda^\top H \lambda - (v + b)^\top \lambda + \frac{1}{2} q^\top Q^\dagger q & \text{for } \lambda \in \Lambda \\ +\infty & \text{otherwise.} \end{cases} \quad (29)$$

A computation like in the proof of Lemma 3.1 shows that

$$\partial\delta(\lambda) = H\lambda - v - b + C(N(Q)), \quad \text{for } \lambda \in \Lambda.$$

Since \mathcal{S}_D is defined as the set of minimizers of δ , one finds

$$\mathcal{S}_D = \{\bar{\lambda} \in \Lambda : H\bar{\lambda} \in v + b + C(N(Q))\}.$$

It is useful to introduce

$$\sigma := \inf_{\substack{\mu \in \partial B \cap R(C) \\ C^\top \mu \in R(Q^\dagger)}} \mu^\top H \mu, \quad (30)$$

which, by definition, takes the value $+\infty$ when $\{\mu \in R(C) : C^\top \mu \in R(Q^\dagger)\} = \{0\}$. Below, the smallest *nonzero* eigenvalue of a *zero* matrix is defined to be $+\infty$.

Lemma 4.1. *The value in $\bar{\mathbb{R}}$ defined by (30) satisfies $\sigma > 0$. It is the smallest nonzero eigenvalue of H when $R(C^\top) \subset R(Q^\dagger)$.*

Proof. We only have to consider the case when $\{\mu \in R(C) : C^\top \mu \in R(Q^\dagger)\} \neq \{0\}$. Then, $Q^\dagger \neq 0$ (because $C^\top \mu = 0$ and $\mu \in R(C)$ imply that $\mu = 0$) and $C \neq 0$. Now, when $C^\top \mu \in R(Q^\dagger) = N(Q^\dagger)^\perp$, $\mu^\top H \mu = \mu^\top C Q^\dagger C^\top \mu \geq \zeta_{\min}(Q^\dagger) \|C^\top \mu\|^2$, where $\zeta_{\min}(Q^\dagger)$ is the smallest nonzero eigenvalue of Q^\dagger . On the other hand, when $\mu \in R(C)$, $\|C^\top \mu\| \geq \sigma_{\min}(C) \|\mu\|$, where $\sigma_{\min}(C)$ is the smallest nonzero singular value of C . We have shown that

$$\sigma \geq \zeta_{\min}(Q^\dagger) \sigma_{\min}(C)^2 > 0.$$

Suppose now that $R(C^\top) \subset R(Q^\dagger)$. Then $\sigma = \inf\{\mu^\top H \mu : \mu \in \partial B \cap R(C)\}$, so that σ will be the smallest nonzero eigenvalue of H if we show that $R(C) = N(H)^\perp$ or that $N(H) = N(C^\top)$. The inclusion $N(C^\top) \subset N(H)$ is clear. Conversely, let $\nu \in N(H)$, which reads $C Q^\dagger C^\top \nu = 0$. This implies that $C^\top \nu \in N(Q^\dagger) = R(Q^\dagger)^\perp \subset N(C)$, by assumption. Then $C^\top \nu = 0$. \square

Note that when $R(C^\top) \not\subset R(Q^\dagger)$, σ is not the smallest nonzero eigenvalue of H . Here is an example

$$Q^\dagger = \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} \quad \text{and} \quad C = \begin{pmatrix} 1 & 1 \end{pmatrix}.$$

Then $H = 1$, while $\sigma = +\infty$ since there is no nonzero μ such that $C^\top \mu \in R(Q^\dagger)$.

Proposition 4.2. *Consider problem (28) with $Q \succcurlyeq 0$ and suppose that it has a solution. Then property (27) is satisfied by the dual function (29), with the Euclidean norm and a constant $L = 1/\sigma$, where σ is defined by (30).*

Proof. Since problem (28) has a solution and its constraint is qualified, \mathcal{S}_D is nonempty (it is identical to the set of optimal multipliers). To prove (27), we only have to consider the dual variables $\lambda \in \Lambda$, since otherwise $\partial\delta(\lambda)$ is empty. Note also that we only have to consider the case when $H \neq 0$ since otherwise $\mathcal{S}_D = \Lambda$ and (27) is trivially satisfied with $L = 0$.

Let $\lambda \in \text{dom } \delta = \Lambda$, $g \in \partial\delta(\lambda)$, and $\bar{\lambda}$ be the projection of λ onto \mathcal{S}_D . We have for some u and $\bar{u} \in N(Q)$

$$g = H\lambda - (v + b) + Cu \quad \text{and} \quad 0 = H\bar{\lambda} - (v + b) + C\bar{u}.$$

Subtracting these two identities and taking the scalar product with $(\lambda - \bar{\lambda})$ yield

$$g^\top(\lambda - \bar{\lambda}) = (\lambda - \bar{\lambda})^\top H(\lambda - \bar{\lambda}) + (u - \bar{u})^\top C^\top(\lambda - \bar{\lambda}).$$

The last term vanishes, since $C^\top(\lambda - \bar{\lambda}) \in R(Q)$ and $Q(u - \bar{u}) = 0$. Now observe that $\lambda - \bar{\lambda} \in R(C)$, since $\bar{\lambda} + N(C^\top) \subset \mathcal{S}_D$ (equality holds actually), and that $C^\top(\lambda - \bar{\lambda}) \in R(Q) = R(Q^\dagger)$, since both λ and $\bar{\lambda} \in \Lambda$. Therefore

$$g^\top(\lambda - \bar{\lambda}) \geq \left(\inf_{\substack{\mu \in \partial B \cap R(C) \\ C^\top \mu \in R(Q^\dagger)}} \mu^\top H \mu \right) \|\lambda - \bar{\lambda}\|^2 = \sigma \|\lambda - \bar{\lambda}\|^2.$$

Now (27) with $L = 1/\sigma$ follows by using the Cauchy-Schwarz inequality on the left hand side. \square

When C is surjective and $Q \succ 0$, extending this result to the dual function associated with the strictly convex quadratic problem (2) is an exercise (then H is positive definite and L is the inverse of the smallest eigenvalue of H). On the other hand, when C is not surjective, property (27) cannot hold without being lightly weakened, as shown by the following example.

Example 4.3. Consider the special QP with a single inactive constraint ($m = 1$, $C = 0$, and $l < 0 < u$) and a zero optimal value ($q = 0$). Then δ is the function

$$\delta(\lambda) = \begin{cases} -u\lambda & \text{if } \lambda \leq 0 \\ -l\lambda & \text{if } \lambda > 0. \end{cases}$$

Clearly, (27) can hold only if λ is not too far from the dual solution set: $|\lambda| \leq L \min(-l, u)$. \square

The analysis of the inequality constrained QP is more difficult than the one of problem (28), since it has to cover simultaneously two different cases: the quadratic dual function of the equality constrained QP (28) and the sharp dual function of the previous example.

In the case of an inequality constrained QP, it will be shown that the Lipschitz constant L may also depend on the largest gap Δ between the $\bar{y} \in \mathcal{S}_P^y$ and the inactive bounds. More precisely, Δ is defined by

$$\Delta := \sup_{\bar{y} \in \mathcal{S}_P^y} \min \left(\min_{i \in I_i \cup J} (u_i - \bar{y}_i), \min_{i \in J \cup I_u} (\bar{y}_i - l_i) \right), \quad (31)$$

where the index sets I_l , J , and I_u are introduced in (13). If $J = \emptyset$, either I_l or $I_u \neq \emptyset$, and the fact that $l < u$ implies that $\Delta > 0$. If $J \neq \emptyset$, the convexity of \mathcal{S}_P implies that there is a $\bar{y} \in \mathcal{S}_P^y$ such that $l_J < \bar{y}_J < u_J$, in which case also $\Delta > 0$. The dependence of L on Δ is clearly visible in Example 4.3: for λ at a unit distance from the solution, we must have $L \geq 1/\min(-l, u) = 1/\Delta$. This lower bound on L goes to infinity when l or u tends to zero, and it goes to zero when $l \rightarrow -\infty$ and $u \rightarrow +\infty$.

Proposition 4.4. *Consider problem (1) with $Q \succcurlyeq 0$ and suppose that it has a solution. Then, for any bounded set $\mathcal{B} \subset \mathbb{R}^m$, there exists a constant L , such that*

$$\forall \lambda \in \mathcal{S}_D + \mathcal{B}, \forall g \in \partial\delta(\lambda) : \text{dist}(\lambda, \mathcal{S}_D) \leq L\|g\|. \quad (32)$$

Proof. *First stage: definition of L .* We know that $\mathcal{S}_D \neq \emptyset$. Let \mathcal{B} be a bounded set in \mathbb{R}^m , i.e., $\mathcal{B} \subset \beta B$ for some $\beta > 0$. To make the proof rigorous, we now define $L > 0$, even though the motivation for its definition will not look quite clear at this point.

Let \mathcal{K} be the collection of index sets $K \subset \{1, \dots, m\}$ such that $I_l \subset K$ and $I_u \subset K^c := \{1, \dots, m\} \setminus K$ (the index sets I_l and I_u are defined in (13)). With any index set $K \subset \{1, \dots, m\}$, we associate the orthant

$$\mathcal{O}_K := \{\lambda \in \mathbb{R}^m : \lambda_K \geq 0, \lambda_{K^c} \leq 0\}.$$

Define \mathcal{O} and \mathcal{A} by (14). For any index set $K \in \mathcal{K}$, $\mathcal{O}_K \cap \mathcal{A}$ is nonempty (since it contains $\mathcal{S}_D = \mathcal{O} \cap \mathcal{A}$, see Lemma 3.2). Therefore, with an index set $K \in \mathcal{K}$, Corollary 3.6 associates two constants $\tau_K > 0$ and $\gamma_K > 0$ such that for any $\lambda \in \mathbb{R}^m$:

$$\text{dist}(\lambda, \mathcal{O}_K) \leq \tau_K \text{dist}(\lambda, \mathcal{O}_K \cap \mathcal{A}) \implies \text{dist}(\lambda, \mathcal{A}) \geq \gamma_K \text{dist}(\lambda, \mathcal{O}_K \cap \mathcal{A}).$$

Since \mathcal{K} is finite, the constants

$$\tau := \min_{K \in \mathcal{K}} \tau_K \quad \text{and} \quad \gamma := \min_{K \in \mathcal{K}} \gamma_K$$

are positive. Therefore, we have found two constants $\tau > 0$ and $\gamma > 0$ such that, for any $K \in \mathcal{K}$, there holds

$$\begin{aligned} \lambda \in \mathcal{O}_K^\tau &:= \{\lambda' \in \mathbb{R}^m : \text{dist}(\lambda', \mathcal{O}_K) \leq \tau \text{dist}(\lambda', \mathcal{O}_K \cap \mathcal{A})\} \\ \implies \text{dist}(\lambda, \mathcal{A}) &\geq \gamma \text{dist}(\lambda, \mathcal{O}_K \cap \mathcal{A}). \end{aligned} \quad (33)$$

Recall the definitions (30) of $\sigma > 0$ and (31) of $\Delta > 0$. Then, the constant $L \geq 0$ is defined by

$$L := \max\left(\frac{1}{\sigma\gamma^2}, \frac{\beta}{\tau\Delta}\right). \quad (34)$$

In this formula, the constants $\sigma > 0$ and $\Delta > 0$ may take the value $+\infty$. Therefore, L is finite, but can vanish.

Second stage: Proof of (32). Fix $\lambda \in \mathcal{S}_D + \mathcal{B}$ and $g \in \partial\delta(\lambda)$ (necessarily, $\lambda \in \text{dom } \delta$). Denote the projection of λ onto \mathcal{S}_D by $\bar{\lambda} := P_{\mathcal{S}_D}(\lambda)$. Observe that,

$$\|\lambda - \bar{\lambda}\| \leq \beta. \quad (35)$$

Let $\varepsilon \in]0, \Delta[$ and define L_ε by formula (34), but with $\Delta - \varepsilon$ in place of Δ . Observe now that showing

$$g^\top(\lambda - \bar{\lambda}) \geq \frac{1}{L_\varepsilon} \|\lambda - \bar{\lambda}\|^2 \quad (36)$$

suffices to conclude the proof since then the inequality in (32) follows from the Cauchy-Schwarz inequality applied to the left hand side of (36) and the fact that ε can be chosen arbitrarily close to zero.

From the form of the subdifferential $\partial\delta(\lambda)$ given by Lemma 3.1, we have for some $y_\lambda \in Y_\lambda$, some $y_{\bar{\lambda}} \in Y_{\bar{\lambda}}$, and some $u, \bar{u} \in N(Q)$:

$$g = H\lambda - v - y_\lambda + Cu \quad \text{and} \quad 0 = H\bar{\lambda} - v - y_{\bar{\lambda}} + C\bar{u}. \quad (37)$$

According to Lemma 3.2, $y_{\bar{\lambda}}$ can be chosen arbitrarily in \mathcal{S}_p^y and we take it such that

$$\min \left(\min_{i \in I_l \cup J} (u_i - \bar{y}_i), \min_{i \in J \cup I_u} (\bar{y}_i - l_i) \right) \geq \Delta - \varepsilon. \quad (38)$$

As in the proof of Proposition 4.2, $(u - \bar{u})^\top C^\top(\lambda - \bar{\lambda}) = 0$, because $C^\top(\lambda - \bar{\lambda}) \in R(Q)$ (both λ and $\bar{\lambda} \in \text{dom } \delta \subset \Lambda$) and $Q(u - \bar{u}) = 0$. Therefore, subtracting the identities in (37) and taking the scalar product with $(\lambda - \bar{\lambda})$ yield

$$g^\top(\lambda - \bar{\lambda}) = (\lambda - \bar{\lambda})^\top H(\lambda - \bar{\lambda}) - (y_\lambda - y_{\bar{\lambda}})^\top(\lambda - \bar{\lambda}). \quad (39)$$

We will get (36) by finding a lower bound of the right hand side of (39). Note that the two terms are nonnegative (this is clear for the first one, since H is positive semi-definite; for the second one, use the monotonicity property (10)).

Since $\bar{\lambda} = P_{\mathcal{A} \cap \mathcal{O}}(\lambda)$, by Lemma 3.7, one can find an index set $K \subset \mathcal{K}$ such that $\bar{\lambda} = P_{\mathcal{A} \cap \mathcal{O}_K}(\lambda)$. We analyze successively two complementary cases, using the set \mathcal{O}_K^τ defined in (33) and $\lambda^t := (1-t)\bar{\lambda} + t\lambda$ for $t \in \mathbb{R}$.

Case A: there exists a $t \in]0, 1]$ such that $\lambda^t \in \mathcal{O}_K^\tau$. In this case, we work on the first term in the right hand side of (39), discarding the second one. Because $P_{\mathcal{A} \cap \mathcal{O}_K}(\lambda^t) = \bar{\lambda}$, (33) gives

$$\gamma \|\lambda^t - \bar{\lambda}\| \leq \|\lambda^t - P_{\mathcal{A}}(\lambda^t)\|.$$

Decompose $\lambda^t - \bar{\lambda} = \mu_0 + \mu_1$, where $\mu_0 \in N(C^\top)$ and $\mu_1 \in R(C)$, and observe that $C^\top \mu_1 = C^\top(\lambda^t - \bar{\lambda}) \in R(Q) = R(Q^\dagger)$ (since both λ^t and $\bar{\lambda} \in \text{dom } \delta \subset \Lambda$). Then, using the definition (30) of σ , one finds

$$(\lambda^t - \bar{\lambda})^\top H(\lambda^t - \bar{\lambda}) = \mu_1^\top H \mu_1 \geq \sigma \|\mu_1\|^2.$$

From the fact that $\mu_1 \in R(C)$ and that $C^\top(\lambda^t - \mu_1) = C^\top(\bar{\lambda} + \mu_0) = C^\top \bar{\lambda} = Q\bar{x} + q$, we deduce that $\mu_1 = \lambda^t - P_{\mathcal{A}}(\lambda^t)$. Therefore

$$(\lambda^t - \bar{\lambda})^\top H(\lambda^t - \bar{\lambda}) \geq \sigma \gamma^2 \|\lambda^t - \bar{\lambda}\|^2.$$

Since $\lambda - \bar{\lambda} = (\lambda^t - \bar{\lambda})/t$, we also have

$$(\lambda - \bar{\lambda})^\top H(\lambda - \bar{\lambda}) \geq \sigma \gamma^2 \|\lambda - \bar{\lambda}\|^2.$$

Discarding the second term in the right hand side of (39) (it is nonnegative) and using the definition of L in (34), it follows that

$$g^\top(\lambda - \bar{\lambda}) \geq (\lambda - \bar{\lambda})^\top H(\lambda - \bar{\lambda}) \geq \sigma\gamma^2 \|\lambda - \bar{\lambda}\|^2 \geq \frac{1}{L} \|\lambda - \bar{\lambda}\|^2 \geq \frac{1}{L_\varepsilon} \|\lambda - \bar{\lambda}\|^2,$$

which is the expected inequality (36).

Case B: for any $t \in]0, 1]$, $\lambda^t \notin \mathcal{O}_K^\tau$. In this case, we work on the second term in the right hand side of (39), discarding the first one. Let us start by choosing $t \in]0, 1]$ sufficiently small such that $\lambda_i^t \bar{\lambda}_i > 0$ when $\bar{\lambda}_i \neq 0$; by assumption, this $\lambda^t \notin \mathcal{O}_K^\tau$. Let $g^t \in \partial\delta(\lambda^t)$, so that $g^t = H\lambda^t - v - y_{\lambda^t} + Cu^t$ for some $y_{\lambda^t} \in Y_{\lambda^t}$ and some $u^t \in N(Q)$ (compare with (37)). Proceeding as before, we get an identity like (39):

$$(g^t)^\top(\lambda^t - \bar{\lambda}) = (\lambda^t - \bar{\lambda})^\top H(\lambda^t - \bar{\lambda}) - (y_{\lambda^t} - y_{\bar{\lambda}})^\top(\lambda^t - \bar{\lambda}). \quad (40)$$

Denote by $\lambda_K^t := P_{\mathcal{O}_K}(\lambda^t)$ the projection of λ^t onto \mathcal{O}_K and decompose

$$-(y_{\lambda^t} - y_{\bar{\lambda}})^\top(\lambda^t - \bar{\lambda}) = -(y_{\lambda^t} - y_{\bar{\lambda}})^\top(\lambda^t - \lambda_K^t) - (y_{\lambda^t} - y_{\bar{\lambda}})^\top(\lambda_K^t - \bar{\lambda}).$$

The last term in the right hand side is nonnegative. Indeed, by the choice of t , if $\bar{\lambda}_i \neq 0$, one has $\lambda_i^t \bar{\lambda}_i > 0$ and therefore $(y_{\lambda^t} - y_{\bar{\lambda}})_i = 0$ (see the definition (9) of Y_λ). The only nonzero terms of the last scalar product are therefore of the form $(y_{\bar{\lambda}} - y_{\lambda^t})_i (\lambda_K^t)_i$. If $(\lambda_K^t)_i > 0$, one has $\lambda_i^t > 0$ (since λ_K^t is the projection of λ^t onto \mathcal{O}_K) and therefore $(y_{\lambda^t})_i = l_i$, so that the term can be written $((y_{\bar{\lambda}})_i - l_i)(\lambda_K^t)_i \geq 0$ (since $l_i \leq (y_{\bar{\lambda}})_i \leq u_i$). Similarly, the term is nonnegative when $(\lambda_K^t)_i < 0$. Therefore

$$-(y_{\lambda^t} - y_{\bar{\lambda}})^\top(\lambda^t - \bar{\lambda}) \geq -(y_{\lambda^t} - y_{\bar{\lambda}})^\top(\lambda^t - \lambda_K^t) = \sum_{i \in I_{K, \lambda^t}} (y_{\bar{\lambda}} - y_{\lambda^t})_i (\lambda^t - \lambda_K^t)_i,$$

where we have introduced the index set

$$I_{K, \lambda^t} := \{i : \lambda_i^t \neq (\lambda_K^t)_i\}.$$

Let us show that all the terms of the sum on the indices $i \in I_{K, \lambda^t}$ above are positive. Observe first that $(\lambda_K^t)_i = 0$ (since $\lambda_i^t \neq (\lambda_K^t)_i$ and λ_K^t is the projection of λ^t onto the orthant \mathcal{O}_K). On the other hand, if $\lambda_i^t > 0$, then $(y_{\lambda^t})_i = l_i$ and $l_i < (y_{\bar{\lambda}})_i \leq u_i$ (this is because $i \in K^c$ when $\lambda_i^t > 0$ and $(\lambda_K^t)_i = 0$, and because $I_l \subset K$); therefore $(y_{\bar{\lambda}} - y_{\lambda^t})_i = (y_{\bar{\lambda}})_i - l_i \geq \Delta - \varepsilon > 0$ (see (38)). Similarly, if $\lambda_i^t < 0$, then $(y_{\lambda^t})_i = u_i$, $i \in K$, and $(y_{\bar{\lambda}} - y_{\lambda^t})_i = (y_{\bar{\lambda}})_i - u_i \leq -(\Delta - \varepsilon) < 0$. In particular, all the terms of the sum are positive. Therefore, we get ($\|\cdot\|_1$ denotes the ℓ_1 -norm)

$$-(y_{\lambda^t} - y_{\bar{\lambda}})^\top(\lambda^t - \bar{\lambda}) \geq (\Delta - \varepsilon) \|\lambda^t - \lambda_K^t\|_1 \geq (\Delta - \varepsilon) \|\lambda^t - \lambda_K^t\| \geq \tau(\Delta - \varepsilon) \|\lambda^t - \bar{\lambda}\|, \quad (41)$$

where the last inequality comes from the fact that $\lambda^t \notin \mathcal{O}_K^\tau$ and that $P_{\mathcal{A} \cap \mathcal{O}_K}(\lambda^t) = \bar{\lambda}$. We

can now conclude:

$$\begin{aligned}
 g^\top(\lambda - \bar{\lambda}) &\geq (g^t)^\top(\lambda - \bar{\lambda}) && \text{[monotonicity of the subdifferential]} \\
 &= \frac{1}{t}(g^t)^\top(\lambda^t - \bar{\lambda}) && \text{[definition of } \lambda^t\text{]} \\
 &\geq -\frac{1}{t}(y_{\lambda^t} - y_{\bar{\lambda}})^\top(\lambda^t - \bar{\lambda}) && \text{[(40)]} \\
 &\geq \frac{\tau(\Delta - \varepsilon)}{t} \|\lambda^t - \bar{\lambda}\| && \text{[(41)]} \\
 &\geq \tau(\Delta - \varepsilon) \|\lambda - \bar{\lambda}\| && \text{[definition of } \lambda^t\text{]} \\
 &\geq \frac{\tau(\Delta - \varepsilon)}{\beta} \|\lambda - \bar{\lambda}\|^2 && \text{[(35)]} \\
 &\geq \frac{1}{L_\varepsilon} \|\lambda - \bar{\lambda}\|^2 && \text{[definition of } L_\varepsilon\text{]}.
 \end{aligned}$$

This is the expected inequality (36). □

4.2. Global linear convergence

We can now state the global linear convergence of the constraint norm towards zero in the AL algorithm of Section 2. Note that the rate of convergence $\min(L/r_k, 1)$ may depend through L on the distance from the initial iterate λ_0 to the dual solution set \mathcal{S}_D .

Theorem 4.5. *Suppose that problem (1) with $Q \succcurlyeq 0$ has a solution. Consider the AL algorithm of Section 2. For any $\beta > 0$, there exists an $L > 0$, such that $\text{dist}(\lambda_0, \mathcal{S}_D) \leq \beta$ implies that*

$$\|y_{k+1} - Cx_{k+1}\| \leq \min\left(\frac{L}{r_k}, 1\right) \|y_k - Cx_k\|, \quad \text{for all } k \geq 1. \quad (42)$$

In particular, if $r_k \geq \bar{r}$ for all $k \geq 1$ and some $\bar{r} > L$, the constraint norm tends to zero globally linearly.

Proof. The proof gathers known techniques (see for example [27, 28]) with the result of Proposition 4.4. We give the details for completeness.

Let us note $g_{k+1} := Cx_{k+1} - y_{k+1}$. Recall from (16) that $g_{k+1} \in \partial\delta(\lambda_{k+1})$. Subtracting two consecutive iteration identities (6) provides

$$\frac{1}{r_{k+1}}(\lambda_{k+2} - \lambda_{k+1}) + (g_{k+2} - g_{k+1}) = \frac{1}{r_k}(\lambda_{k+1} - \lambda_k).$$

Taking norms, using the monotonicity of the subdifferential (which implies that $(g_{k+2} - g_{k+1})^\top(\lambda_{k+2} - \lambda_{k+1}) \geq 0$), and discarding $\|g_{k+2} - g_{k+1}\|^2 \geq 0$, we get $\|\lambda_{k+2} - \lambda_{k+1}\|^2/r_{k+1}^2 \leq \|\lambda_{k+1} - \lambda_k\|^2/r_k^2$ or $\|g_{k+2}\| \leq \|g_{k+1}\|$. This yields the second part of the min in (42).

Subtracting an arbitrary dual solution $\bar{\lambda} \in \mathcal{S}_D$ from both sides of the iteration identity (6) gives

$$\lambda_{k+1} - \bar{\lambda} + r_k g_{k+1} = \lambda_k - \bar{\lambda}, \quad \text{for } k \geq 0.$$

Taking norms and using the monotonicity of the subdifferential lead to

$$\|\lambda_{k+1} - \bar{\lambda}\|^2 + r_k^2 \|g_{k+1}\|^2 \leq \|\lambda_k - \bar{\lambda}\|^2, \quad \text{for } k \geq 0. \quad (43)$$

This shows in particular that the sequence $\{\|\lambda_k - \bar{\lambda}\|\}_{k \geq 0}$ is nonincreasing. Since $\bar{\lambda}$ is arbitrary in \mathcal{S}_D , there holds

$$\text{dist}(\lambda_{k+1}, \mathcal{S}_D) \leq \|\lambda_{k+1} - P_{\mathcal{S}_D}(\lambda_k)\| \leq \|\lambda_k - P_{\mathcal{S}_D}(\lambda_k)\| = \text{dist}(\lambda_k, \mathcal{S}_D).$$

Therefore, $\{\text{dist}(\lambda_k, \mathcal{S}_D)\}_{k \geq 0}$ is also nonincreasing, so that $\lambda_k \in \mathcal{S}_D + \beta B$ for all $k \geq 0$. Now, let $L > 0$ be the constant that Proposition 4.4 associates with $\mathcal{B} := \beta B$. By this proposition, $\|\lambda_k - P_{\mathcal{S}_D}(\lambda_k)\| \leq L \|g_k\|$. Discarding the first term in the left hand side of (43) and using $P_{\mathcal{S}_D}(\lambda_k)$ for $\bar{\lambda}$, we get $\|g_{k+1}\| \leq (L/r_k) \|g_k\|$. This yields the first part of the min in (42). \square

5. Numerical experiments and discussion

The aim of this section is to illustrate by numerical experiments the global linear convergence property of the AL algorithm studied in this paper and to assess the quality of the bound given by Theorem 4.5. The numerical experiments are taken from seismic reflection tomography applications. We conclude with a discussion on algorithmic implications.

5.1. A seismic reflection tomography problem

Seismic reflection tomography is a technique used to recover the geological structure of the subsoil from the measurements of the travel-times of seismic waves (see [10] for a description of the approach). From an optimization viewpoint, the problem consists in minimizing a nonlinear least-squares function subject to nonlinear constraints. In [5], a Gauss-Newton SQP method globalized by line-search is proposed and analyzed. At each iteration, a solution to a strictly convex quadratic model of the objective function subject to linear constraints is computed using our code QPAL [4, 6].

We have chosen here to present the results obtained with the problem KARINE, which is representative of those observed with our collection of 2D and 3D seismic reflection problems. These have a number of variables up to $15 \cdot 10^3$ and a number of constraints up to 10^4 . The features of the selected problem are summarized in Table 5.1. It is a 2D model

n	m	m_{act}^*	κ_2
442	320	108	$8.4 \cdot 10^5$

Table 5.1: Description of the tomography problem KARINE

depending on $n = 442$ parameters and having $m = 320$ linear inequality constraints. The matrix Q of the selected quadratic subproblem (1) is positive definite and has its ℓ_2 condition number equal to $\kappa_2 = 8.4 \cdot 10^5$. Its constraint matrix C has been balanced (the Euclidean norm of each of its rows is equal to 1). The number of active constraints at the solution is $m_{act}^* = 108$, which represents 34 % of the number of constraints.

The results presented in Section 5.2 have been obtained using the AL algorithm described in Section 2, with a fixed augmentation parameter r . In order to study the dependence of the results on r , we have run the QP solver for 21 different values of r , ranging from 1 to 10^5 . In each case, the AL algorithm is initialized with a null Lagrangian multiplier ($\lambda_0 = 0$) and is stopped when the constraint norm is sufficiently small ($\|y_k - Cx_k\| \leq 10^{-10}$).

5.2. Assessing the global linear convergence result

In this section, we illustrate the global linear convergence property of the AL algorithm established in Theorem 4.5. The actual global rate of linear convergence is given by $\rho := \sup\{\|y_{k+1} - Cx_{k+1}\|/\|y_k - Cx_k\| : k \geq 1\}$ and can be estimated during a particular run by

$$\rho_{\text{est}} := \max_{1 \leq k \leq n_{AL}} \frac{\|y_{k+1} - Cx_{k+1}\|}{\|y_k - Cx_k\|} \leq \rho, \tag{44}$$

where n_{AL} is the number of AL iterations actually performed to reach the required accuracy of 10^{-10} on the constraint norm.

Theorem 4.5 has shown that ρ is bounded above by a function of r :

$$\rho \leq \min\left(\frac{L}{r}, 1\right) \quad \text{or} \quad \log \rho \leq \min(\log L - \log r, 0). \tag{45}$$

A natural question is to know whether this bound is tight in practice. This is difficult to say, since the value of L is generally unknown, but the appearance of ρ_{est} as a function of r in the considered problem may give a clue on this question.

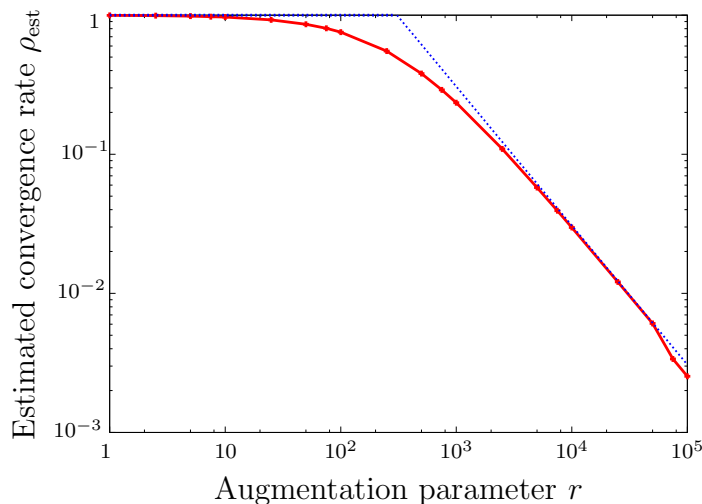


Figure 5.1: Global linear convergence rate of the constraint norm as a function of r

The plain curve in Figure 5.1 gives $\log \rho_{\text{est}}$ as a function of $\log r$ (double logarithmic scale). As predicted by the theory, we see that $\rho_{\text{est}} \leq 1$ for all positive r . Furthermore, the larger is the augmentation parameter r , the faster is the convergence: $\rho \simeq 1$ for $r \leq 10$ (convergence is hardly detectable) and $\rho \simeq 3 \cdot 10^{-3}$ for $r = 10^5$ (convergence is obtained in very few AL iterations). We have represented by a dotted line the tangent to the ρ_{est} curve with a slope -1 . This line crosses the top horizontal line of the graph at the horizontal coordinate $L_{\text{inf}} \simeq 304$. According to (45), L_{inf} provides a lower estimate of the value of L . Since both curves (the plain and dotted ones) are quite close, it is likely that the dotted curve is close to the upper bound on ρ given by (45). As a result, it is likely that the upper bound given by Theorem 4.5 is tight. Note that the small discrepancy between both curves for large values of r ($r \geq 7.5 \cdot 10^4$) is due here to the fact that the AL algorithm reaches the required constraint norm accuracy in very few AL iterations ($n_{AL} \leq 3$), so that the maximum in (44) is taken on that few number. In other cases,

such a discrepancy can come from an inexact solve of the bound constraint problem (5), due to a large value of r .

5.3. Discussion

As shown in this paper, the global linear rate of convergence of the AL algorithm depends on the Lipschitz constant L given by (34) and on the value of $\bar{r} := \inf r_k$, where r_k is the value of the augmentation parameter at iteration k . More precisely, the decrease of the constraint norm at iteration k is bounded above by L/r_k . It is usually impossible to compute L in practice, since it depends on the constants γ , σ , and Δ (see Lemma 3.4, (30), (31), and finally (34)), which are either unknown or too expensive to compute. As a Lipschitz constant, however, L has easily computable *lower* estimates.

The estimate L_{inf} of L given in Section 5.2 is not available at run time, since it requires to run the AL algorithm on a particular problem for various values of r . Nevertheless, the quantities

$$L_{\text{inf},k} := \max_{1 \leq i \leq k} \left(r_i \frac{\|y_{i+1} - Cx_{i+1}\|}{\|y_i - Cx_i\|} \right)$$

satisfy $L_{\text{inf},k} \leq L$ and can therefore be used as a lower estimate of L , after iteration k is completed. A given desired rate of convergence $\rho_{\text{des}} \in]0, 1[$ is then likely to be obtained at iteration $k + 1$ by taking

$$r_{k+1} \geq \frac{L_{\text{inf},k}}{\rho_{\text{des}}}. \quad (46)$$

It is the fact that the estimate (42) has a *global* validity that gives sense to an update of the value of r_k in this way at *each* iteration. It should be clear at this point that the AL algorithm gains in efficiency by taking r_k as large as possible, the only limitation being that problem (5) needs to be numerically solvable. Since it is sometimes difficult to tell what is a large value for a particular problem, the lower bound on r_{k+1} in (46) may also be useful as a reference.

We conclude with a result providing an estimate of the number of iterations needed to reach a given tolerance on the constraint norm. Assume that a number $\rho_{\text{des}} \in]0, 1[$ is given as a desired rate of convergence. Of course, since the Lipschitz constant L is unknown, this rate of convergence cannot be ensured, but the algorithm can try to approach it by updating r_k when it feels it is necessary. The next result gives an estimate of the iterative complexity of the AL algorithm with an update rule based on (46). More precisely, defining

$$\rho_k := \frac{\|y_{k+1} - Cx_{k+1}\|}{\|y_k - Cx_k\|},$$

the AL algorithm is supposed to update the value of r_k , for $k \geq 1$, according to:

$$\text{if } \rho_k \leq \rho_{\text{des}}, \quad \text{then } r_{k+1} = r_k, \quad \text{else } r_{k+1} = \frac{\rho_k}{\rho_{\text{des}}} r_k. \quad (47)$$

There is nothing magic in the update rule of r_k above. It could equally use $r_{k+1} = 10 \rho_k r_k / \rho_{\text{des}}$ or simply $r_{k+1} = 10 r_k$ when r_k needs to be increased.

Proposition 5.1. *Suppose that the AL algorithm of Section 2 uses the rule (47) to update the augmentation parameter r_k , for $k \geq 1$. Let $\varepsilon \in]0, 1[$ and let L be the positive constant*

given by Theorem 4.5. Fix any $t \in]\rho_{\text{des}}, 1[$. Then

$$\|y_{k+1} - Cx_{k+1}\| \leq \varepsilon \|y_1 - Cx_1\|, \quad (48)$$

as soon as

$$k \geq \frac{\log \varepsilon}{\log t} + \max \left(1 + \frac{\log(L/(tr_1))}{\log(t/\rho_{\text{des}})}, 0 \right). \quad (49)$$

Proof. Let $t \in]\rho_{\text{des}}, 1[$. Clearly, since $\rho_i \leq 1$,

$$\frac{\|y_{k+1} - Cx_{k+1}\|}{\|y_1 - Cx_1\|} = \prod_{1 \leq i \leq k} \rho_i \leq \prod_{\substack{1 \leq i \leq k \\ \rho_i \leq t}} \rho_i \leq t^{k_t},$$

where $k_t := |K_t|$ is the number of elements in $K_t := \{i \in \mathbb{N} : 1 \leq i \leq k, \rho_i \leq t\}$. Taking logarithms, we see that (48) holds as soon as $k_t \geq (\log \varepsilon)/(\log t)$.

If $K_t^c := \{1, \dots, k\} \setminus K_t$ is empty, then $k = k_t$ and the result is proven.

Suppose now that $K_t^c \neq \emptyset$. Since $\rho_i \leq L/r_i$ (by Theorem 4.5), $i \in K_t$ as soon as $r_i \geq L/t$. Let j be the last index in K_t^c , namely the $(k - k_t)$ th one, if any. Then r_j is the result of $k - k_t - 1$ updates from r_1 , using factors ρ_i/ρ_{des} that are $\geq t/\rho_{\text{des}}$ (see the update rule (47)). Hence we must have $(t/\rho_{\text{des}})^{k - k_t - 1} r_1 \leq r_j \leq L/t$. This gives an upper bound on the number of elements of K_t^c , namely

$$k - k_t \leq 1 + \frac{\log(L/(tr_1))}{\log(t/\rho_{\text{des}})}.$$

The total number of iterations to satisfy (48) is therefore at most this upper bound on the number of elements in K_t^c , plus the lower bound on k_t obtained above. \square

Roughly expressed, the number of iterations needed to reach precision $\varepsilon > 0$ on the relative constraint norm is of order $O(\log \varepsilon) + O(\log L)$. As shown in the proof of Proposition 5.1, the first term of order $O(\log \varepsilon)$ is due to the linear convergence of the constraint norm towards zero, which is triggered when the augmentation parameter is large enough (a consequence of Theorem 4.5). The second term of order $O(\log L)$, which is the only place where the dimension of the problem can intervene, is due to a possible too small value of r_1 and to the number of iterations that the rule (47) needs to make r_k large enough. This term can be made as small as desired by choosing a large value for r_1 or by adopting an update rule of r_k that increases these values more rapidly than in (47). As a result, the *computational* complexity of the AL algorithm of Section 2 essentially rests on the one of the AL subproblems (5). When strict complementarity holds, the finite identification of the active constraints in (5) occurs and the computational complexity is then basically induced by the very first AL subproblems. Our experience with the AL algorithm, limited to the seismic reflection tomography problems described in Section 5.1, supports that conclusion.

References

- [1] K. J. Arrow, R. M. Solow: Gradient methods for constrained maxima with weakened assumptions, in: Studies in Linear and Nonlinear Programming, K. J. Arrow et al. (eds.), Stanford University Press, Standford (1958) 166–176.

- [2] J.F. Bonnans, J. Ch. Gilbert, C. Lemaréchal, C. Sagastizábal: Numerical Optimization – Theoretical and Practical Aspects, Springer, Berlin (2003).
- [3] J.V. Burke, M.C. Ferris: Characterization of solution sets of convex programs, *Oper. Res. Lett.* 10 (1991) 57–60.
- [4] F. Delbos: Problèmes d’Optimisation de Grande Taille avec Contraintes en Tomographie de Réflexion 3D, PhD Thesis, University Pierre et Marie Curie (Paris VI), Paris, 2004.
- [5] F. Delbos, J. Ch. Gilbert, R. Glowinski, D. Sinoquet: Constrained optimization in seismic reflection tomography: An SQP augmented Lagrangian approach, *Geophys. J. Int.* (2005), submitted.
- [6] F. Delbos, J. Ch. Gilbert, D. Sinoquet: QPAL: A solver of large-scale convex quadratic optimization problems using an augmented Lagrangian approach, Technical report, INRIA, Le Chesnay, 2005, to appear.
- [7] D. den Hertog: Interior Point Approach to Linear, Quadratic and Convex Programming, *Mathematics and its Applications 277*, Kluwer Academic Publishers, Dordrecht (1992).
- [8] M. Fortin, R. Glowinski: Méthodes de Lagrangien Augmenté – Applications à la Résolution Numérique de Problèmes aux Limites, *Méthodes Mathématiques de l’Informatique 9*, Dunod, Paris (1982).
- [9] A. Friedlander, J.M. Martínez: On the maximization of a concave quadratic function with box constraints, *SIAM J. Optim.* 4 (1994) 177–192.
- [10] R. Glowinski, Q.-H. Tran: Constrained optimization in reflection tomography: The augmented Lagrangian method, *East-West J. Numer. Math.* 1(3) (1993) 213–234.
- [11] M.R. Hestenes: Multiplier and gradient methods, *J. Optimization Theory Appl.* 4 (1969) 303–320.
- [12] J.-B. Hiriart-Urruty, C. Lemaréchal: Convex Analysis and Minimization Algorithms, *Grundlehren der mathematischen Wissenschaften 305-306*, Springer (1993).
- [13] A.F. Izmailov, M.V. Solodov: Error bounds for 2-regular mappings with Lipschitzian derivatives and their applications, *Math. Program.* 89 (2001) 413–435.
- [14] B. Jansen: Interior Point Techniques in Optimization – Complementarity, Sensitivity and Algorithms, *Applied Optimization 6*, Kluwer Academic Publishers, Dordrecht (1997).
- [15] M. Kojima, N. Megiddo, T. Noma, A. Yoshise: A Unified Approach to Interior Point Algorithms for Linear Complementarity Problems, *Lecture Notes in Computer Science 538*, Springer, Berlin (1991).
- [16] O.L. Mangasarian: A simple characterization of solution sets of convex programs, *Oper. Res. Lett.* 7 (1988) 21–26.
- [17] J.J. Moré, G. Toraldo: On the solution of large quadratic programming problems with bound constraints, *SIAM J. Optim.* 1 (1991) 93–113.
- [18] Y.E. Nesterov, A.S. Nemirovskii: Interior-Point Polynomial Algorithms in Convex Programming, *SIAM Studies in Applied Mathematics 13*, SIAM, Philadelphia (1994).
- [19] J. Nocedal, S.J. Wright: Numerical Optimization, Springer Series in Operations Research, Springer, New York (1999).
- [20] J.-S. Pang: Methods for quadratic programming: A survey, *Computers and Chemical Engineering* 7 (1983) 583–594.
- [21] J.-S. Pang: Error bounds in mathematical programming, *Math. Program.* 79 (1997) 299–332.

- [22] M. J. D. Powell: A method for nonlinear constraints in minimization problems, in: Optimization, R. Fletcher (ed.), Academic Press, London (1969) 283–298.
- [23] R. T. Rockafellar: Convex Analysis, Princeton University Press, Princeton (1970).
- [24] R. T. Rockafellar: New applications of duality in convex programming, in: Proceedings of the 4th Conference of Probability, Brasov, Romania (1971) 73–81.
- [25] R. T. Rockafellar: A dual approach to solving nonlinear programming problems by unconstrained optimization, *Math. Program.* 5 (1973) 354–373.
- [26] R. T. Rockafellar: Augmented Lagrange multiplier functions and duality in nonconvex programming, *SIAM J. Control* 12 (1974) 268–285.
- [27] R. T. Rockafellar: Monotone operators and the proximal point algorithm, *SIAM J. Control Optimization* 14 (1976) 877–898.
- [28] R. T. Rockafellar: Augmented Lagrangians and applications of the proximal point algorithm in convex programming, *Math. Oper. Res.* 1 (1976) 97–116.
- [29] C. Roos, T. Terlaky, J.-Ph. Vial: Theory and Algorithms for Linear Optimization – An Interior Point Approach, John Wiley & Sons, Chichester (1997)
- [30] T. Terlaky (ed.): Interior Point Methods of Mathematical Programming, Kluwer Academic Press, Dordrecht (1996).
- [31] S. J. Wright: Primal-Dual Interior-Point Methods, SIAM, Philadelphia (1997).