

Sequential equality-constrained optimization for nonlinear programming*

E. G. Birgin[†] L. F. Bueno[‡] J. M. Martínez[§]

September 9, 2015[¶]

Abstract

A novel idea is proposed for solving optimization problems with equality constraints and bounds on the variables. In the spirit of Sequential Quadratic Programming and Sequential Linearly-Constrained Programming, the new proposed approach approximately solves, at each iteration, an equality-constrained optimization problem. The bound constraints are handled in outer iterations by means of an Augmented Lagrangian scheme. Global convergence of the method follows from well-established nonlinear programming theories. Numerical experiments are presented.

Key words: Nonlinear programming, Sequential Equality-Constrained Optimization, Augmented Lagrangian, numerical experiments.

1 Introduction

Although nonlinearly constrained optimization is a well-established area of numerical mathematics, many challenges remain and stimulate the development of new methods.

*This work was supported by PRONEX-CNPq/FAPERJ E-26/111.449/2010-APQ1, FAPESP (grants 2010/10133-0, 2013/03447-6, 2013/05475-7, 2013/07375-0, and 2015/02528-8), and CNPq (grants 309517/2014-1 and 303750/2014-6).

[†]Department of Computer Science, Institute of Mathematics and Statistics, University of São Paulo, São Paulo, SP, Brazil. e-mail: egbirgin@ime.usp.br

[‡]Institute of Science and Technology, Federal University of São Paulo, São José dos Campos, SP, Brazil. e-mail: lfelipebueno@gmail.com

[§]Department of Applied Mathematics, Institute of Mathematics, Statistics, and Scientific Computing, State University of Campinas, Campinas, SP, Brazil. e-mail: martinez@ime.unicamp.br

[¶]Revisions made on February 4 and April 19, 2016.

Frequently, novel algorithms are firstly developed for equality-constrained optimization and, later, their extension to equality and inequality constraints is introduced. Handling inequality constraints may lead to cumbersome combinatorial problems, if one decides to rely on active-set strategies. In Sequential Quadratic Programming (SQP) algorithms [16, 29], a quadratic function is minimized on a polytope at each iteration. As a consequence, the behavior of SQP algorithms is strongly affected by the performance of the Quadratic Programming solver. In Sequential Linearly-Constrained Programming [28], one minimizes a (non-necessarily quadratic) Lagrangian (or Augmented Lagrangian) on the polytope that represents the linearization of the constraints. Again, this is a considerably complicated subproblem in which all the combinatorial issues are incorporated to handle inequalities. In the most popular Augmented Lagrangian (AL) methods [1, 9, 13] equality constraints (and sometimes also inequality constraints [1]) are incorporated into the objective function of the subproblems, so that the only inequality constraints of subproblems are represented by the bounds on the variables, thus simplifying the combinatorial difficulties. Finally, in some Interior-Point methods [32], subproblems involve the minimization of a barrier function that tends to infinity on the boundaries, subject to the equality constraints. In this case, the behavior near the boundary is problematic. In Linear Programming, long experience with Interior-Point methods taught software developers how to deal with this inconvenience switching to the “central trajectory”, but this is not so simple when we deal with nonlinear constraints.

On the other hand, minimization with only equality constraints is a very attractive subproblem. The main reason is that the standard optimality condition for this problem is a nonlinear system of equations, instead of the system with equalities and inequalities that appear in the KKT conditions of general optimization. In many cases, solving equality constrained subproblems is less difficult than solving quadratic programming subproblems with inequalities of similar dimensions.

This state of facts led us to define a method for general nonlinear programming with equality constrained subproblems, in which the bound constraints on the variables are incorporated into the objective function under the Augmented Lagrangian interaction. At each outer iteration the method minimizes an equality-constrained optimization problem, using a mild stopping criterion. For this task we use Newton-like linearizations of the Lagrangian nonlinear system. After the outer iteration, the Lagrange multipliers and penalty parameters corresponding to bound constraints are checked and updated as in the AL scheme described in [9]. This is the first implementation of an Augmented Lagrangian method in which the subproblem constraints (or lower-level constraints in [1]) adopt a nontrivial definition, which, by the way, reinforce the necessity of the general theory of [1, 9].

We describe an implementation of these ideas in an algorithm called SECO (Sequential Equality-Constrained Optimization). Section 2 describes the main algorithm. The method for equality-constrained nonlinear minimization that is used for solving the sub-

problems is described in Section 3. Numerical experiments are described and analyzed in Section 4. Conclusions are given in Section 5.

Notation. If $x \in \mathbb{R}^n$, $x_+ = \max\{x, 0\}$ where the maximum is taken componentwise.

2 SECO algorithm

In this work we address the problem

$$\text{Minimize } f(x) \text{ subject to } h(x) = 0 \text{ and } \ell \leq x \leq u,$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $h : \mathbb{R}^n \rightarrow \mathbb{R}^m$ are sufficiently smooth and $\ell \leq x \leq u$ must hold componentwise, i.e. $\ell_i \leq x_i \leq u_i$, $i = 1, \dots, n$, with $\ell = (\ell_1, \dots, \ell_n)^T$ and $u = (u_1, \dots, u_n)^T$. (Note that, paying the price of adding slack variables, inequality constraints of the form $g(x) \leq 0$ can be incorporated as $g(x) + s = 0$ plus $s \geq 0$.) However, to simplify the presentation and since the extension is straightforward, in the rest of this work we will deal with the problem

$$\text{Minimize } f(x) \text{ subject to } h(x) = 0 \text{ and } x \geq 0. \quad (1)$$

For all $x \in \mathbb{R}^n$, $\rho > 0$, and $v \in \mathbb{R}_+^n$, we define the Augmented Lagrangian function $L_\rho(x, v)$ (that only penalizes the non-negativity constraints) by

$$L_\rho(x, v) = f(x) + \frac{\rho}{2} \sum_{v_i/\rho \geq x_i} (v_i/\rho - x_i)^2.$$

Therefore, we have that

$$\nabla L_\rho(x, v) = \nabla f(x) - \sum_{v_i/\rho \geq x_i} (v_i - \rho x_i) e^i$$

and

$$\nabla^2 L_\rho(x, v) = \nabla^2 f(x) + \rho \sum_{v_i/\rho \geq x_i} e^i (e^i)^T,$$

where $e^i \in \mathbb{R}^n$ stands for the i -th canonical vector in \mathbb{R}^n . The operators ∇ and ∇^2 always indicate derivatives with respect to x . The Hessian $\nabla^2 L_\rho(x, v)$ is discontinuous at the points in which $x_i = v_i/\rho$ but we can adopt the definition above without leading to contradictions.

The main model algorithm that we now describe is a particular case of Algorithm 4.1 of [9, p.33].

Algorithm 2.1: SECO

Let $v_{\max} > 0$, $\gamma > 1$, $0 < \tau < 1$, $\bar{v}^1 \in [0, v_{\max}]^n$, and $\rho_1 > 0$ be given. Initialize $k \leftarrow 1$.

Step 1. Find $x^k \in \mathbb{R}^n$ as an approximate solution of

$$\text{Minimize } L_{\rho_k}(x, \bar{v}^k) \text{ subject to } h(x) = 0. \quad (2)$$

Step 2. Compute $V^k = \min\{x^k, \bar{v}^k/\rho_k\}$. If $k = 1$ or

$$\|V^k\|_{\infty} \leq \tau \|V^{k-1}\|_{\infty}, \quad (3)$$

choose $\rho_{k+1} \geq \rho_k$. Otherwise, choose $\rho_{k+1} \geq \gamma\rho_k$.

Step 3. Choose $\bar{v}^{k+1} \in [0, v_{\max}]^n$.

Step 4. Set $k \leftarrow k + 1$ and go to Step 1.

Remark 1. At Step 1, one tries to find an approximate solution of (2). In fact, at iteration k , we require a point $x^k \in \mathbb{R}^n$ such that there exists $\lambda^k \in \mathbb{R}^m$ satisfying

$$\begin{aligned} \|\nabla L_{\rho_k}(x^k, \bar{v}^k) + \nabla h(x^k)\lambda^k\|_{\infty} &\leq \varepsilon_k^{\text{opt}} \\ \|h(x^k)\|_{\infty} &\leq \varepsilon_k^{\text{feas}}, \end{aligned} \quad (4)$$

where $\{\varepsilon_k^{\text{feas}}\}$ and $\{\varepsilon_k^{\text{opt}}\}$ are sequences that tend to zero in a way that will be specified later.

Remark 2. In theory, the only requirement for the multipliers \bar{v}^{k+1} is that they belong to a given compact interval. However, in order to speed up convergence, in practice, we define

$$v^{k+1} = (\bar{v}^k - \rho_k x^k)_+ \quad (5)$$

and compute \bar{v}^{k+1} as the projection of v^{k+1} onto the safeguarding interval $[0, v_{\max}]^n$. Moreover, it is worthwhile to mention that, in the theorem below, $\{\lambda^k\}$ in (4) and $\{v^{k+1}\}$ in (5) play the role of Lagrange multipliers associated with equality and inequality constraints, respectively.

Definition 2.1. [5] Assume that the constraints of an optimization problem are $h(x) = 0$ and $g(x) \leq 0$, x^* is a feasible point, and I is the set of indices of active inequality constraints at x^* . For all $x \in \mathbb{R}^n$ we define $K(x)$ as the non-negative cone generated by $\{\pm \nabla h_i(x), \forall i\}$ and $\{\nabla g_i(x), i \in I\}$. The Cone-Continuity property is said to hold at x^* when the point-to-set mapping $K(x)$ is continuous at x^* .

Theorem 2.1. Assume that $\{x^k\}$ is a sequence generated by SECO. Then,

1. Every limit point of $\{x^k\}$ satisfies the AKKT (Approximate Karush-Kuhn-Tucker) condition [2] of the problem

$$\text{Minimize } \sum_{i=1}^n \max\{0, -x_i\}^2 \text{ subject to } h(x) = 0. \quad (6)$$

2. If a limit point of $\{x^k\}$ satisfies the Cone-Continuity property with respect to $h(x) = 0$, then this limit point satisfies the KKT conditions of the problem (6).
3. If a limit point of $\{x^k\}$ is feasible for problem (1), then it satisfies the AKKT conditions of (1).
4. If a limit point of $\{x^k\}$ is feasible and satisfies the Cone-Continuity property with respect to the constraints $h(x) = 0, x \geq 0$, then this point satisfies the KKT conditions of (1).

Proof. Parts 1 and 3 of the thesis follow from Theorems 6.3 and 6.2 of [9], respectively; while parts 2 and 4 follow from Theorem 3.3 of [9]. It is worth noting that Theorem 3.3 of [9] uses the U-condition; while parts 2 and 4 in the thesis use the Cone-Continuity property. The U-condition and the Cone-Continuity property are equivalent. The Cone-Continuity property is the geometrical interpretation of the U-condition (see [5] for details). \square

The main result of [5] is that the Cone-Continuity property is the weakest constraint qualification that guarantees that AKKT implies KKT. Therefore, it is weaker than constraint qualifications LICQ (Linear Independence of the gradients of active constraints), Mangasarian-Fromovitz [25], CPLD (Constant Positive Linear Dependence) [31, 6], RCPLD (Relaxed Positive Linear Dependence) [3], and CPG (Constant Positive Generators) [4].

By the AKKT property and Theorem 6.1 of [9], given $\varepsilon^{\text{feas}} > 0$ and $\varepsilon^{\text{opt}} > 0$, there exists k such that, on success, $(x^k, \lambda^k, v^{k+1})$ satisfies

$$\begin{aligned}
\|\nabla f(x^k) - \sum_{i=1}^n v^{k+1} e^i + \nabla h(x^k) \lambda^k\|_\infty &\leq \varepsilon^{\text{opt}}, \\
\|h(x^k)\|_\infty &\leq \varepsilon^{\text{feas}}, \\
-x^k &\leq \varepsilon^{\text{feas}}, \\
\min\{x^k, v^{k+1}\} &\leq \varepsilon^{\text{feas}}.
\end{aligned} \tag{7}$$

3 Algorithm for equality-constrained minimization

In this section, we describe the method being proposed to approximately solve the subproblems (2) of the SECO algorithm.

In order to avoid cumbersome notation, we denote $F(x) = L_{\rho_k}(x, \bar{v}^k)$. Then, F has Lipschitz-continuous first derivatives and $\nabla F(x)$ is semismooth [30]. Problem (2) becomes

$$\text{Minimize } F(x) \text{ subject to } h(x) = 0. \tag{8}$$

Moreover, for the algorithm described in this section we will also use the index k to identify iterations and iterates, which should not be confused with the ones of the main algorithm described in the previous section. We now describe a method that, basically, consists of solving problem (8) by applying Newton's method to its KKT conditions.

The KKT conditions of problem (8) are given by

$$\begin{aligned}\nabla\mathcal{L}(x, \lambda) &= 0 \\ h(x) &= 0,\end{aligned}\tag{9}$$

where $\mathcal{L}(x, \lambda) = F(x) + h(x)^T\lambda$, and, thus, $\nabla\mathcal{L}(x, \lambda) = \nabla F(x) + \nabla h(x)\lambda$.

The Newtonian linear system associated with the nonlinear system of equations (9) is given by

$$\begin{pmatrix} \nabla^2\mathcal{L}(x, \lambda) & \nabla h(x) \\ \nabla h(x)^T & 0 \end{pmatrix} \begin{pmatrix} d_x \\ d_\lambda \end{pmatrix} = - \begin{pmatrix} \nabla\mathcal{L}(x, \lambda) \\ h(x) \end{pmatrix}.\tag{10}$$

However, the linear system that we solve at each iteration is of the form

$$\begin{pmatrix} \nabla^2\mathcal{L}(x, \lambda) + \epsilon^{\text{nw}}I & \nabla h(x) \\ \nabla h(x)^T & -\epsilon^{\text{se}}I \end{pmatrix} \begin{pmatrix} d_x \\ d_\lambda \end{pmatrix} = - \begin{pmatrix} \nabla\mathcal{L}(x, \lambda) \\ h(x) \end{pmatrix},\tag{11}$$

where $\epsilon^{\text{nw}}, \epsilon^{\text{se}} \geq 0$ are real numbers such that the coefficient matrix in (11) has n positive eigenvalues and m negative eigenvalues (and no null eigenvalues). In principle, no modification of the Newtonian linear system (10) would be necessary. However, we are interested in minimizers of (8), and not on other kind of solutions of the Lagrange conditions (9) such as maximizers or other stationary points. In order to increase the chance of convergence to minimizers (or, at least, discourage convergence to other stationary points), the matrix of the system is modified in such a way that the modified Hessian of the Lagrangian $\nabla^2\mathcal{L}(x, \lambda) + \epsilon^{\text{nw}}I$ is positive definite onto the null space of $\nabla h(x)^T$. This goal is achieved with the modification displayed in (11) of the diagonal of the coefficient matrix in (10), which corresponds to the modification of its inertia. On the other hand, when $\epsilon^{\text{se}} > 0$, the diagonal matrix $-\epsilon^{\text{se}}I$ ensures that the last m rows of the coefficient matrix in (11) are linearly independent. Moreover, if $\epsilon^{\text{se}} > 0$, the Sylvester Law of Inertia [18, p. 403] and the identity

$$\begin{pmatrix} \nabla^2\mathcal{L}(x, \lambda) + \epsilon^{\text{nw}}I + \nabla h(x)\nabla h(x)^T/\epsilon^{\text{se}} & 0 \\ 0 & -\epsilon^{\text{se}}I \end{pmatrix} = \begin{pmatrix} I & \nabla h(x)/\epsilon^{\text{se}} \\ 0 & I \end{pmatrix} \begin{pmatrix} \nabla^2\mathcal{L}(x, \lambda) + \epsilon^{\text{nw}}I & \nabla h(x) \\ \nabla h(x)^T & -\epsilon^{\text{se}}I \end{pmatrix} \begin{pmatrix} I & 0 \\ \nabla h(x)^T/\epsilon^{\text{se}} & I \end{pmatrix},\tag{12}$$

imply that the matrix $\nabla^2\mathcal{L}(x, \lambda) + \epsilon^{\text{nw}}I + \nabla h(x)\nabla h(x)^T/\epsilon^{\text{se}}$ is positive definite. As it will be shown later, this implies that the algorithm being proposed is well defined.

At iteration k , given the iterate (x^k, λ^k) , a linear system like (11) is solved to find (d_x^k, d_λ^k) . Ideally, we would like to compute the new iterate (x^{k+1}, λ^{k+1}) as $(x^k, \lambda^k) + (d_x^k, d_\lambda^k)$. However, in practice, the next iterate will be given by

$$\begin{pmatrix} x^{k+1} \\ \lambda^{k+1} \end{pmatrix} = \begin{pmatrix} x^k \\ \lambda^k \end{pmatrix} + \alpha_k \begin{pmatrix} \gamma_x^k d_x^k \\ \gamma_\lambda^k d_\lambda^k \end{pmatrix},$$

where $\gamma_x^k, \gamma_\lambda^k \in (0, 1]$ have the role of limiting the size of the search directions $\tilde{d}_x^k = \gamma_x^k d_x^k$ and $\tilde{d}_\lambda^k = \gamma_\lambda^k d_\lambda^k$, respectively, and $\alpha_k \in (0, 1]$ is computed in order to obtain simple decrease of a merit function not yet specified.

Therefore, the three main ingredients of the developed Newton's method applied to the KKT system (9) of problem (8) are: (a) the correction of the inertia of the coefficient matrix in (16); (b) the control step (computation of γ_x^k and γ_λ^k); and (c) the backtracking procedure to compute α_k and the decision of the merit function considered in the backtracking.

The description of the algorithm requires the definition of the squared infeasibility measurement given by

$$\Phi(x) = \frac{1}{2} \|h(x)\|_2^2. \quad (13)$$

We also define, in a rather unconventional way, the Augmented Lagrangian function associated with problem (8) as

$$L_\epsilon(x, \lambda) = \epsilon \mathcal{L}(x, \lambda) + \Phi(x), \quad (14)$$

where $\epsilon \geq 0$ is given. The whole method is described in the algorithm below.

Algorithm 3.1.

Let $\varepsilon^{\text{feas}}, \varepsilon^{\text{opt}} > 0$ be given tolerances related to the stopping criteria. Let \hat{x}^0 and λ^0 be initial estimates of the primal and dual variables, respectively. Let `PERTURBX0` $\in \{\text{TRUE}, \text{FALSE}\}$ and `NOBCKTRCKATALL` $\in \{\text{TRUE}, \text{FALSE}\}$ be given parameters. Set $k \leftarrow 0$, `NOMORESTPCNTRL` $\leftarrow \text{FALSE}$, and `NOMOREBCKTRCK` $\leftarrow \text{FALSE}$.

Step 0. *Perturbation (or not) of initial guess*

If `PERTURBX0` then set

$$x_i^0 = \hat{x}_i^0 + 0.01 \xi_i |\hat{x}_i^0|, \text{ for } i = 1, \dots, n,$$

where ξ_i is a random variable with uniform distribution within the interval $[-1, 1]$ for $i = 1, \dots, n$. Otherwise, set $x^0 = \hat{x}^0$.

Step 1. *Stopping criteria*

Step 1.1 If

$$\|h(x^k)\|_\infty \leq \varepsilon^{\text{feas}} \text{ and } \|\nabla \mathcal{L}(x^k, \lambda^k)\|_\infty \leq \varepsilon^{\text{opt}},$$

stop by declaring that an approximate KKT point with the required feasibility and optimality tolerances has been found.

Step 1.2 If

$$\|h(x^k)\|_\infty \leq \epsilon^{\text{feas}} \text{ and } f(x^k) \leq -10^{10},$$

stop by declaring that the objective function appears to be unbounded from below.

Step 2. Inertia correction

Step 2.1.

Case $k = 0$: Set $\epsilon^{\text{nw}} \leftarrow 0$. If $m \leq n$, set $\epsilon^{\text{se}} \leftarrow 0$, otherwise set $\epsilon^{\text{se}} \leftarrow 10^{-8}$.

Case $k \geq 1$: If $\gamma_x^{k-1} = 1$, set $\epsilon^{\text{nw}} \leftarrow 0.1\epsilon_{k-1}^{\text{nw}}$, otherwise set $\epsilon^{\text{nw}} \leftarrow 3\epsilon_{k-1}^{\text{nw}}$. Set $\epsilon^{\text{se}} \leftarrow 0.1\epsilon_{k-1}^{\text{se}}$.

Step 2.2. Compute the inertia (i_+, i_-, i_0) , where i_+ , i_- , and i_0 represent the number of positive, negative, and null eigenvalues, respectively, of the matrix

$$\begin{pmatrix} \nabla^2 \mathcal{L}(x^k, \lambda^k) + \epsilon^{\text{nw}} I & \nabla h(x^k) \\ \nabla h(x^k)^T & -\epsilon^{\text{se}} I \end{pmatrix}. \quad (15)$$

Step 2.3. If $i_+ < n$ then increase ϵ^{nw} by setting $\epsilon^{\text{nw}} \leftarrow \max\{10^{-8}, 3\epsilon^{\text{nw}}\}$.

Step 2.4. If $i_- < m$ then increase ϵ^{se} by setting $\epsilon^{\text{se}} \leftarrow \max\{10^{-8}, 3\epsilon^{\text{se}}\}$.

Step 2.5. If ϵ^{nw} or ϵ^{se} were increased at Steps 2.3 or 2.4, respectively, go to Step 2.2.

Step 2.6. Define $\epsilon_k^{\text{nw}} = \epsilon^{\text{nw}}$ and $\epsilon_k^{\text{se}} = \epsilon^{\text{se}}$.

Step 3. Search direction and step control

Step 3.1. Compute d_x^k and d_λ^k as the (unique) solution of the linear system

$$\begin{pmatrix} \nabla^2 \mathcal{L}(x^k, \lambda^k) + \epsilon_k^{\text{nw}} I & \nabla h(x^k) \\ \nabla h(x^k)^T & -\epsilon_k^{\text{se}} I \end{pmatrix} \begin{pmatrix} d_x^k \\ d_\lambda^k \end{pmatrix} = - \begin{pmatrix} \nabla \mathcal{L}(x^k, \lambda^k) \\ h(x^k) \end{pmatrix}. \quad (16)$$

Step 3.2. If `NO MORE STP CNTRL`, assign $\gamma_x^k = 1$ and $\gamma_\lambda^k = 1$. Otherwise, define

$$\gamma_x^k = \min \left\{ 1, \frac{100 \max\{1, \|x^k\|_\infty\}}{\|d_x^k\|_\infty} \right\} \text{ and } \gamma_\lambda^k = \min \left\{ 1, \frac{100 \max\{1, \|\lambda^k\|_\infty\}}{\|d_\lambda^k\|_\infty} \right\}.$$

Step 3.3. Define $\tilde{d}_x^k = \gamma_x^k d_x^k$ and $\tilde{d}_\lambda^k = \gamma_\lambda^k d_\lambda^k$.

Step 4. Backtracking (on primal variables only)

Step 4.1. If `NO BCKTRCK AT ALL` or `NO MORE BCKTRCK`, set $\alpha_k = 1$ and go to Step 5.

Step 4.2. Set $\alpha \leftarrow 1$.

Step 4.3. Set $x_{\text{trial}} \leftarrow x^k + \alpha \tilde{d}_x^k$.

Step 4.4. Consider conditions

$$\|h(x^k)\|_\infty \leq \varepsilon^{\text{feas}} \text{ and } \epsilon_k^{\text{se}} = 0 \text{ and } \text{fl}(\mathcal{L}(x_{\text{trial}}, \lambda^k)) \leq \text{fl}(\mathcal{L}(x^k, \lambda^k)) \quad (17)$$

and

$$(\|h(x^k)\|_\infty > \varepsilon^{\text{feas}} \text{ or } \epsilon_k^{\text{se}} > 0) \text{ and } \text{fl}(L_{\epsilon_k^{\text{se}}}(x_{\text{trial}}, \lambda^k)) \leq \text{fl}(L_{\epsilon_k^{\text{se}}}(x^k, \lambda^k)), \quad (18)$$

where $\text{fl}(\cdot)$ represents the result of performing an operation in floating point arithmetic.

If (17) or (18) hold, define $\alpha_k = \alpha$ and go to Step 5.

Step 4.5. Set $\alpha \leftarrow \alpha/2$ and go to Step 4.3.

Step 5. *Update the iterate*

Step 5.1. Define $x^{k+1} = x^k + \alpha_k \tilde{d}_x^k$ and $\lambda^{k+1} = \lambda^k + \alpha_k \tilde{d}_\lambda^k$.

Step 5.2. Consider conditions

$$\text{fl}(\mathcal{L}(x^{k+1}, \lambda^k)) = \text{fl}(\mathcal{L}(x^k, \lambda^k)), \quad (19)$$

and

$$\text{fl}(L_{\epsilon_k^{\text{se}}}(x^{k+1}, \lambda^k)) = \text{fl}(L_{\epsilon_k^{\text{se}}}(x^k, \lambda^k)). \quad (20)$$

If (17,19) or (18,20) hold, set $\text{NOMORESTPCNTRL} \leftarrow \text{TRUE}$ and $\text{NOMOREBCKTRCK} \leftarrow \text{TRUE}$.

Step 5.3. Set $k \leftarrow k + 1$ and go to Step 1.

Remarks. It is well-known that backtracking may be harmful in some well characterized situations. This is why the algorithm includes the possibility of avoiding backtracking by setting the logical parameter `NOBCKTRCKATALL`. Perturbation of the initial guess may be adequate to avoid some undesirable situations in which symmetry prevents convergence. This is the reason why the method includes the possibility of perturbing the initial guess by setting the logical parameter `PERTURBX0`.

In Algorithm 3.1 (Step 4.4) the acceptance of the trial step (see (17) and (18)) is conditioned to the decrease of two different merit functions: the Lagrangian $\mathcal{L}(\cdot, \lambda^k)$ and the Augmented Lagrangian $L_{\epsilon_k^{\text{se}}}(\cdot, \lambda^k)$. Consider the three cases **(i)** $\epsilon_k^{\text{se}} = 0$ and $h(x^k) = 0$, **(ii)** $\epsilon_k^{\text{se}} = 0$ and $h(x^k) \neq 0$, and **(iii)** $\epsilon_k^{\text{se}} > 0$. We will prove that, in case (i), the direction d_x^k is a descent direction for the Lagrangian; while, in cases (ii) and (iii), the direction d_x^k is a descent direction for the Augmented Lagrangian.

In case (i), the second block of equations in (16) with $\epsilon_k^{\text{se}} = 0$ and $h(x^k) = 0$, says that d_x^k belongs to the null-space of $\nabla h(x^k)^T$. Therefore, pre-multiplying the first block of equations in (16) by $(d_x^k)^T$ and using, that, due to the inertia correction, $\nabla^2 \mathcal{L}(x^k, \lambda^k) + \epsilon_k^{\text{se}} I$ is positive definite on the null-space of $\nabla h(x^k)^T$, we obtain that $(d_x^k)^T \nabla \mathcal{L}(x^k, \lambda^k)$ is negative. Thus, when $\epsilon_k^{\text{se}} = 0$ and $h(x^k) = 0$, d_x^k is a descent direction for the Lagrangian $\mathcal{L}(\cdot, \lambda^k)$ at x^k .

In case (ii), since $\epsilon_k^{\text{se}} = 0$, the Augmented Lagrangian $L_{\epsilon_k^{\text{se}}}(\cdot, \lambda^k)$ coincides with $\Phi(\cdot)$. Premultiplying the equation $\nabla h(x^k)^T d_x^k + h(x^k) = 0$ by $h(x^k)$, since $h(x^k) \neq 0$, we obtain that the scalar product of the direction with the $\nabla\Phi(x^k)$ (that coincides with $\nabla L_{\epsilon_k^{\text{se}}}(x^k, \lambda^k)$) is negative, as desired.

In case (iii), by the second block of equations in (16), we have that $\nabla h(x^k)^T d_x^k - \epsilon_k^{\text{se}} d_\lambda^k = -h(x^k)$. Thus, $d_\lambda^k = (\nabla(x^k)^T d_x^k + h(x^k))/\epsilon_k^{\text{se}}$. Now, by the first block of equations in (16),

$$[\nabla^2 \mathcal{L}(x^k, \lambda^k) + \epsilon_k^{\text{nw}} I] d_x^k + \nabla h(x^k) d_\lambda^k = -\nabla \mathcal{L}(x^k, \lambda^k).$$

Therefore,

$$[\nabla^2 \mathcal{L}(x^k, \lambda^k) + \epsilon_k^{\text{nw}} I] d_x^k + [\nabla h(x^k) \nabla h(x^k)^T d_x^k + \nabla h(x^k) h(x^k)] / \epsilon_k^{\text{se}} = -\nabla \mathcal{L}(x^k, \lambda^k).$$

Thus,

$$\begin{aligned} [\nabla^2 \mathcal{L}(x^k, \lambda^k) + \epsilon_k^{\text{nw}} I + \nabla h(x^k) \nabla h(x^k)^T / \epsilon_k^{\text{se}}] d_x^k &= -\nabla \mathcal{L}(x^k, \lambda^k) - \nabla \Phi(x^k) / \epsilon_k^{\text{se}} \\ &= -\nabla L_{\epsilon_k^{\text{se}}}(x^k, \lambda^k) / \epsilon_k^{\text{se}}. \end{aligned}$$

But, by the inertia correction procedure, the identity (12), and the Sylvester Law of Inertia, we have that $\nabla^2 \mathcal{L}(x^k, \lambda^k) + \epsilon_k^{\text{nw}} I + \nabla h(x^k)^T \nabla h(x^k) / \epsilon_k^{\text{se}}$ is positive definite. Therefore d_x^k is a descent direction for the Augmented Lagrangian $L_{\epsilon_k^{\text{se}}}(\cdot, \lambda^k)$ at x^k .

Summing up, we have proved that

$$h(x^k) = 0 \text{ and } \epsilon_k^{\text{se}} = 0 \text{ and } \mathcal{L}(x_{\text{trial}}, \lambda^k) \leq \mathcal{L}(x^k, \lambda^k) \quad (21)$$

or

$$(h(x^k) \neq 0 \text{ or } \epsilon_k^{\text{se}} > 0) \text{ and } L_{\epsilon_k^{\text{se}}}(x_{\text{trial}}, \lambda^k) \leq L_{\epsilon_k^{\text{se}}}(x^k, \lambda^k) \quad (22)$$

eventually holds. Conditions (17) and (18) in Algorithm 3.1 are the heuristic ‘‘practical’’ versions of (21) and (22), respectively.

Algorithm 3.1 may be considered as an heuristic Newton-based procedure for solving the equality-constrained minimization problem (8). For obtaining global convergence, this algorithm may be naturally coupled with the globally convergent Flexible Inexact-Restoration (FIR) procedure for solving (8) introduced in [11]. Each iteration of the FIR algorithm, as other Inexact-Restoration methods, has two phases: Feasibility and Optimality. In the Feasibility Phase the algorithm aims to improve feasibility and in the Optimality Phase the algorithm minimizes approximately a Lagrangian approximation subject to the linearization of the constraints [7, 15, 17, 19, 21, 22, 23, 24, 26, 27]. As a consequence, a trial point is obtained which is accepted as new iterate or not according to the value of a Sharp Lagrangian merit function [15, 26]. When the trial point is not accepted, a new trial point is obtained that satisfies the linearized constraints. If the feasible phase is well-defined, it can be proved that the algorithm generates AKKT

sequences and that limit points are KKT when the Cone-Continuity property holds. The conditions required in the Feasibility Phase are mild and thus many different (perhaps heuristic) methods may be employed for this phase. Our proposal here essentially consists of employing Algorithm 3.1 as the Feasibility Phase of the FIR method. More precisely, the globalization procedure can be sketched in the following Algorithm 3.2.

Algorithm 3.2

Initialize the standard algorithmic parameters of FIR and the parameters of Algorithm 3.1. Let $N > 0$, $\zeta \in (0, 1)$, and $v > 0$. Initialize $k \leftarrow 0$.

Step 1. Run Algorithm 3.1 employing a maximum of N iterations and obtaining the point y^{k+1} . If the standard stopping criterion of Algorithm 3.1 is satisfied, stop. Otherwise, if $\|h(y^{k+1})\| \leq \zeta \|h(x^k)\|$ and $\|y^{k+1} - x^k\| \leq v \|h(x^k)\|$ go to Step 2, else stop by declaring Failure in the Feasibility Phase.

Step 2. Proceed as in FIR by updating the penalty parameter, minimizing the Lagrangian in the tangent set and making the necessary comparisons by means of which, eventually, x^{k+1} is obtained. Update $k \leftarrow k + 1$ and go to Step 1.

4 Numerical experiments

We implemented Algorithms 2.1 and 3.1 in Fortran 90. This means that the SECO subproblems will be solved by an heuristic method and that the globalization scheme suggested in the previous section will not be considered; the main reason for this choice has been keeping the implemented method simple. All tests were conducted on a computer with 3.5 GHz Intel Core i7 processor and 16GB 1600 MHz DDR3 RAM memory, running OS X Yosemite (version 10.10.4). Codes were compiled by the GFortran compiler of GCC (version 4.9.2) with the -O3 optimization directive enabled.

Regarding the parameters of Algorithm 2.1, arbitrarily but based on previous experimentation with the Augmented Lagrangian solver Algencan [1, 9], we set $v_{\max} = 10^{20}$, $\gamma = 10$, and $\tau = 0.5$. We also set $\varepsilon^{\text{feas}} = \varepsilon^{\text{opt}} = 10^{-8}$. The choice of the sequences $\{\varepsilon_k^{\text{feas}}\}$ and $\{\varepsilon_k^{\text{opt}}\}$ (see (4)) is given by

$$\varepsilon_1^{\text{feas}} = \sqrt{\varepsilon^{\text{feas}}} \quad \text{and} \quad \varepsilon_1^{\text{opt}} = \sqrt{\varepsilon^{\text{opt}}}$$

and, for $k > 1$, if

$$\|h(x^{k-1})\|_{\infty} \leq \sqrt{\varepsilon^{\text{feas}}} \quad \text{and} \quad \|\nabla \mathcal{L}(x^{k-1}, \lambda^{k-1})\|_{\infty} \leq \sqrt{\varepsilon^{\text{opt}}}$$

then

$$\begin{aligned} \varepsilon_k^{\text{feas}} &= \max\{\varepsilon^{\text{feas}}, \min\{\frac{1}{10}\varepsilon_{k-1}^{\text{feas}}, \frac{1}{2}\|h(x^{k-1})\|_{\infty}\}\} \\ \varepsilon_k^{\text{opt}} &= \max\{\varepsilon^{\text{opt}}, \min\{\frac{1}{10}\varepsilon_{k-1}^{\text{opt}}, \frac{1}{2}\|\nabla \mathcal{L}(x^{k-1}, \lambda^{k-1})\|_{\infty}\}\}. \end{aligned}$$

Otherwise, $\varepsilon_k^{\text{feas}} = \varepsilon_{k-1}^{\text{feas}}$ and $\varepsilon_k^{\text{opt}} = \varepsilon_{k-1}^{\text{opt}}$. The value of the initial penalty parameter ρ_1 , follows exactly the settings considered in Algencan (see [9, p.153]). In this way, differences between SECO and Algencan are concentrated in the choice of lower- and upper-level constraints, making their comparison useful for the purposes of the present work. The stopping criterion of Algorithm 2.1 is given by (7).

Subroutine MA57 from HSL [33] was used to compute the inertia-revealing factorizations¹ at Step 2.2 of Algorithm 3.1 and to solve the linear systems at Step 3.1. The values of parameters `PERTURBX0` \in `{TRUE, FALSE}` and `NOBCKTRCKATALL` \in `{TRUE, FALSE}` will be the subject of numerical experimentation.

4.1 When the initial point should be perturbed

Consider the problem

$$\text{Minimize } (x + y - 10)^2 \text{ subject to } xy = 1. \quad (23)$$

Let $\delta = 5 + 2\sqrt{6}$. It is easy to see that points of the form $(\delta, \delta^{-1})^T$ and $(\delta^{-1}, \delta)^T$ are global minimizers of problem (23). They both annihilate the objective function, satisfy the LICQ constraint qualification (since any feasible point satisfies it), and, hence, satisfy the KKT condition (both with null multiplier). Therefore, it is natural to assume that they should be the target of any optimization method.

Since the point $(5, 5)^T$ is an unconstrained global minimizer (annihilates the objective function but does not satisfy the constraint), it might be a natural initial guess for any iterative solver trying to solve problem (23). However, starting from $\hat{x}^0 = (5, 5)^T$ (with $\lambda^0 = 0$, `PERTURBX0` = `NOBCKTRCKATALL` = `FALSE`), in 9 iterations and never abandoning the line $x = y$, Algorithm 3.1 converges to the KKT point $(1, 1)^T$, that is a local *maximizer* (there is a local minimizer at $(-1, -1)^T$ and the objective function goes to infinity within the feasible set when, for example, x goes to infinity and $y = 1/x$).

Some readers may think that having found a local maximizer is not an issue at all since a KKT point was found. However, from the authors' point of view, this is, at least, an undesired situation, since the Optimizers' main goal is, in most cases, to obtain the lowest possible functional values within the feasible region.

In this subsection we claim that a very simple and affordable way of avoiding the behaviour described above is to slightly perturb the initial guess. When the initial point $\hat{x}^0 = (5, 5)^T$

¹The inertia is freely computed as an outcome of the matrix factorization performed by the Harwell subroutine MA57. Given an indefinite sparse symmetric matrix A , subroutine MA57 computes the factorization given by $PAP^T = LDL^T$, where P is a permutation matrix, L is lower triangular, and D is block diagonal with 1×1 or 2×2 diagonal blocks. By the Sylvester Law of Inertia, matrices A and D have the same inertia and computing the inertia of D is trivial (it is necessary to check the sign or to compute the eigenvalues of its 1×1 or 2×2 diagonal blocks).

is perturbed, as described at Step 0 of Algorithm 3.1 (and the actual initial point given by $x^0 \approx (5.0466, 4.9629)^T$), Algorithm 3.1 (with $\lambda^0 = 0$, PERTURBX0 = TRUE, and NOBACKTRACKATALL = FALSE) converges to the approximate global minimizer $(9.98990, 0.10102)^T$ in 4 iterations.

4.2 When backtracking should be avoided

Algorithm 3.1 presented in Section 3 aims to solve problems of the form (8) with

$$F(x) = L_{\rho_k}(x, \bar{v}^k) = f(x) + \frac{\rho_k}{2} \sum_{v_i/\rho_k \geq x_i} (\bar{v}_i^k/\rho_k - x_i)^2.$$

As a consequence, it is natural to evaluate Algorithm 3.1 considering objective functions of this form in which ρ_k is possibly large. Thus, in the present section, we consider problems of the form

$$\begin{aligned} \text{Minimize} \quad & \frac{1}{2}x^T A x + b^T x + r \sum_{i=1}^n \max\{0, -x_i\}^2 \\ \text{subject to} \quad & \frac{1}{2}x^T A_j x + b_j^T x + c_j = 0, \quad j = 1, \dots, m, \end{aligned} \quad (24)$$

where

$$A = \left[\frac{1}{\max\{1, \max_{i,j}\{\bar{a}_{ij}\}\}} \right] (\bar{A}^T \bar{A}) \text{ and } A_j = \frac{1}{2}(\bar{A}_j^T + \bar{A}_j), \quad (25)$$

and $\bar{A} \in \mathbb{R}^{n \times n}$, $b \in \mathbb{R}^n$, $\bar{A}_j \in \mathbb{R}^{n \times n}$, and $b_j \in \mathbb{R}^n$ for $j = 1, \dots, m$ have random elements with uniform distribution within the interval $[-1, 1]$. In order to generate a feasible problem, a point \bar{x} with random non-negative elements \bar{x}_i with uniform distribution in $[0, 1]$ is generated and $c_j \in \mathbb{R}$ is defined as

$$c_j = -[\frac{1}{2}\bar{x}^T A_j \bar{x} + b_j^T \bar{x}] \text{ for } j = 1, \dots, m. \quad (26)$$

We consider twelve instances of problem (24) with $n = 1000$, $m = 500$, and $r \in \{0, 1, 10, 100, \dots, 10^{10}\}$. The initial point x^0 is always given by

$$x_i^0 = \bar{x}_i + 10^{-8}\xi_i|\bar{x}_i|, \text{ for } i = 1, \dots, n,$$

where ξ_i is a random variable with uniform distribution in $[-1, 1]$, i.e. x^0 is a very slight perturbation of the known feasible point used to generate the problem. At x^0 we have that $\max_{1 \leq j \leq m}\{|x^T A_j x + b_j^T x + c_j|\} \approx 10^{-6} \not\leq \varepsilon^{\text{feas}} = 10^{-8}$. With these settings, the situation we are trying to mimic is the following: (a) Since the Jacobian of the constraints $\nabla h(x^k)$ has full row rank with high probability, we will have $\varepsilon_k^{\text{se}}$ (the perturbation applied to the south-east block of the Jacobian matrix (15) of the KKT system to correct its inertia) equal to 0 for all k ; (b) As a consequence, the primal direction d_x^k will satisfy $h'(x^k)d_x^k + h(x^k) = 0$ for all k and, since $\|h(x^k)\|$ is relatively “small”, we will have $h'(x^k)d_x^k \approx 0$; and (c) With

high probability, the merit function (that coincides with the feasibility measure $\frac{1}{2}\|h(x)\|_2^2$ if x^k is infeasible (i.e. $\|h(x^k)\|_\infty \not\leq \varepsilon^{\text{feas}}$) and $\epsilon_k^{\text{se}} = 0$) will decrease only for very small steps along d_x^k . This *gedanken* experiment illustrates a well-known inconvenience of the usage of merit functions that gave rise to the development of alternative merit functions, nonmonotone strategies, and filter methods, among others. Our claim at the present moment is: try Newton as pure as possible avoiding other alternatives.

Table 1 shows the results of applying Algorithm 3.1 with `PERTURBX0 = FALSE` and `NOBCKTRCKATALL = TRUE` to the twelve instances described in the paragraph above. In the table, `#it` and `#fcnt` stand for number of iterations and number of objective function evaluations, respectively; while “CPU time (s)” is the elapsed CPU time in seconds. In the last two columns, $f(x^*)$ is the objective function value at the last iterate; while $\min_{\{1 \leq i \leq n\}} \{x_i^*\}$ is the smallest component of the final iterate x^* . The results for the case `PERTURBX0 = TRUE` and `NOBCKTRCKATALL = TRUE` are very similar to those displayed in Table 1, i.e. as expected, a small perturbation in the initial guess has no meaningful effect in the behavior of the method, since symmetry is not an issue in this problem.

| r | #it | #fcnt | CPU time (s) | $f(x^*)$ | $\min_{\{1 \leq i \leq n\}} \{x_i^*\}$ |
|-----------|-----|-------|--------------|------------|--|
| 0 | 45 | 46 | 142.92 | 5.5997D+03 | -2.0293D+00 |
| 1 | 74 | 75 | 233.22 | 5.5683D+03 | -2.2747D+00 |
| 10 | 57 | 58 | 179.92 | 7.2960D+03 | -1.6530D+00 |
| 10^2 | 57 | 58 | 174.40 | 1.3483D+04 | -8.7815D-01 |
| 10^3 | 85 | 86 | 268.49 | 1.9558D+04 | -2.0632D-01 |
| 10^4 | 90 | 91 | 273.76 | 2.1172D+04 | -2.9567D-02 |
| 10^5 | 121 | 122 | 377.82 | 2.1376D+04 | -3.1172D-03 |
| 10^6 | 76 | 77 | 239.00 | 2.1418D+04 | -3.2517D-04 |
| 10^7 | 127 | 128 | 397.06 | 2.1656D+04 | -2.8549D-05 |
| 10^8 | 130 | 131 | 412.17 | 2.1953D+04 | -2.5745D-06 |
| 10^9 | 332 | 333 | 1057.63 | 2.1486D+04 | -2.8450D-07 |
| 10^{10} | 169 | 170 | 534.56 | 2.1922D+04 | -2.9008D-08 |

Table 1: Performance of Algorithm 3.1 with `PERTURBX0 = FALSE` and `NOBCKTRCKATALL = TRUE` applied to twelve instances of problem (24) with $r \in \{0, 1, 10, \dots, 10^{10}\}$.

On the other hand, the behavior of the method with `NOBCKTRCKATALL = FALSE` and `PERTURBX0 = FALSE` (in order to preserve the properties of the picked initial guess) is completely different. For the case $r = 0$ (the only one we evaluate), the method gets virtually stuck at the initial guess. In its first 30 iterations, it performs 24 functional evaluations per iteration and takes a step of size $2^{-23} \approx 10^{-7}$. The squared Euclidean norm of the constraints is approximately 1.5275×10^{-10} (the sup-norm is approx. 2.2899×10^{-6})

at the initial guess and, after 30 iterations, those values are exactly the same, showing that the method made no progress at all.

The situation described in the paragraph above represents an extreme case. However, similar results can be observed when the initial guess x^0 has random elements with uniform distribution within the interval $[0, 1]$. With no backtracking at all, Algorithm 3.1, when applied to problem (24) with $r = 0$, satisfies the stopping criterion using 83 iterations (and 84 functional evaluations) using 255.00 seconds of CPU time. Similarly to the values reported in Table 1 for a different initial guess, the objective functional value at the final point is approximately 5.2558×10^3 and the value of the most negative entrance of the final iterate is approximately -2.0929 . On the other hand, with backtracking, the algorithm makes an excruciatingly slow progress satisfying the stopping criterion after 6 792 iterations that require 106 795 functional evaluations (in average, approximately 16 per iterations) and consumes 49 977.22 seconds of CPU time (approx. 14 hours). The achieved objective functional value is approximately 5.7927×10^3 and the value of the most negative entrance of the final iterate is approximately -2.2919 .

4.3 Problems with a variable number of bound constraints

In this section, we consider the problem

$$\begin{aligned} & \text{Minimize} && \frac{1}{2}x^T A x + b^T x \\ & \text{subject to} && \frac{1}{2}x^T A_j x + b_j^T x + c_j = 0, \quad j = 1, \dots, m, \\ & && x_i \geq 0, \quad i \in I, \end{aligned} \tag{27}$$

where $I \subseteq \{1, \dots, n\}$. As in the previous section, A and A_j ($j = 1, \dots, m$) are given by (25) and $\bar{A} \in \mathbb{R}^{n \times n}$, $b \in \mathbb{R}^n$, $\bar{A}_j \in \mathbb{R}^{n \times n}$, and $b_j \in \mathbb{R}^n$ for $j = 1, \dots, m$ have random elements with uniform distribution within the interval $[-1, 1]$. Moreover, in order to obtain a feasible problem, $c_j \in \mathbb{R}$ is given by (26), where \bar{x} has non-negative elements \bar{x}_i with uniform distribution in $[0, 1]$. The considered initial guess x^0 has random elements with uniform distribution within the interval $[0, 1]$.

With the aim of generating a set of instances with solutions with an increasing number of active bound constraints, we consider instances with an increasing number of bound constraints. At a given instance, the probability of $i \in I$ is given by a constant $p \in [0, 1]$ (for $i = 1, \dots, n$) and six instances with $p \in \{0, 0.01, 0.1, 0.5, 0.75, 1\}$ are considered. In all cases we set $n = 1\,000$ and $m = 500$. All instances were solved with Algencon 3.0.0 (inhibiting its acceleration process [8]) and SECO. Both methods found an approximate stationary point in all the six instances. Table 2 shows the results. In the table, “#act” is the number of active bound constraints at the final iterate x^* . In the case of Algencon, this number is given by the cardinality of the set $\{i \mid x_i^* = 0\}$; while in SECO it is given by

the cardinality of the set $\{i \mid |x_i^*| \leq \varepsilon^{\text{feas}} = 10^{-8}\}$ ². Column “#it” displays the number of outer iterations and the number of inner iterations separated by a slash. The next three columns “#fcnt”, “#gcnt”, and “#hcnt” represent the number of functional evaluations and first- (gradients) and second-order (Hessians) derivatives evaluations, respectively. “CPU time (s)” means the elapsed CPU time in seconds and $f(x^*)$ is the value of the objective function at the final iterate. The sup-norm of the equality constraints as well as the sup-norm of the gradient of the Lagrangian are always smaller than $\varepsilon^{\text{feas}}$ and ε^{opt} , respectively (recall that both tolerances are equal to 10^{-8}). Last but not least, bound constraints are satisfied with zero tolerance in Algencan (meaning that $x^* \geq 0$ holds); while SECO satisfies the bound constraints with tolerance $\varepsilon^{\text{feas}}$ (meaning that $x_i^* \geq -\varepsilon^{\text{feas}}$ for $i = 1, \dots, n$). Analyzing the figures in the table, it can be seen that SECO found smaller functional values in the cases $p = 0$ and $p = 1$; while Algencan found smaller functional values in the remaining four cases. If, on the one hand, the smaller functional values found by SECO may be attributed to the admissible enlarged bound constraints, on the other hand Algencan found smaller functional values more times. This means that, in this experiment, Algencan appears to exhibit a larger tendency to find better stationary points than SECO. On the efficiency side, both methods appear to require similar efforts for the case $p = 1$ (instance with the largest number of active constraints at the final iterates); while SECO is considerably faster in the other instances. Thus, disregarding the functional values, if the goal were to find stationary points, using SECO would be a reasonable choice.

4.4 Massive comparisons

4.4.1 Problems with only equality constraints

In this set of experiments, we consider *all* the 190 problems from the CUTEst collection [20] (with their default dimensions and default primal \hat{x}^0 and dual λ^0 initial guesses) that have only equality constraints. Considering this set of problems, we evaluate the performance of Algorithm 3.1 with the four possible combinations of NOBCKTRCKATALL and PERTURBX0 $\in \{\text{TRUE}, \text{FALSE}\}$. A CPU time limit of 5 minutes was imposed.

We used performance profiles [14] to compare the methods. Consider q methods M_1, \dots, M_q and p problems P_1, \dots, P_p and let t_{ij} be a metric of the effort that method M_i made in

²When analyzing the output of Algencan, the reason that motivates using a strict criterion that considers that the primal active variables are those exactly equal to zero is that the subproblems solved by Algencan preserve non-negativity of the primal variables all along the calculations. Moreover, the searches used in the Algencan subproblems’ solver make it almost impossible the existence of positive primal variables with very small values. On the other hand, in SECO, the bound constraints may be violated (hopefully slightly) or variables may be strictly positive with tiny values, so that tolerances are necessary to declare almost feasibility and activity.

| Algencon 3.0.0 without the acceleration process | | | | | | | | |
|---|-------|------|-----------|-------|-------|-------|--------------|------------|
| p | $ I $ | #act | #it | #fcnt | #gcnt | #hcnt | CPU time (s) | $f(x^*)$ |
| 0 | 0 | 0 | 12 / 1083 | 4258 | 1141 | 1083 | 3741.28 | 5.3520D+03 |
| 0.01 | 10 | 3 | 14 / 1175 | 4709 | 1255 | 1172 | 8306.38 | 5.3294D+03 |
| 0.1 | 109 | 36 | 13 / 1185 | 3266 | 1243 | 1175 | 5930.13 | 6.9063D+03 |
| 0.5 | 492 | 182 | 5 / 822 | 1794 | 851 | 812 | 2539.07 | 1.1931D+04 |
| 0.75 | 735 | 253 | 6 / 895 | 1773 | 918 | 883 | 2249.93 | 1.6419D+04 |
| 1 | 1000 | 302 | 7 / 535 | 1034 | 560 | 523 | 1087.35 | 2.1492D+04 |

| SECO (Algorithms 2.1 and 3.1) with no backtracking | | | | | | | | |
|--|-------|------|---------|-------|-------|-------|--------------|------------|
| p | $ I $ | #act | #it | #fcnt | #gcnt | #hcnt | CPU time (s) | $f(x^*)$ |
| 0 | 0 | 0 | 1 / 83 | 84 | 84 | 83 | 261.78 | 5.2558D+03 |
| 0.01 | 10 | 5 | 2 / 66 | 71 | 71 | 66 | 202.94 | 5.7933D+03 |
| 0.1 | 109 | 42 | 3 / 106 | 113 | 113 | 106 | 314.55 | 7.1996D+03 |
| 0.5 | 492 | 183 | 3 / 50 | 57 | 57 | 50 | 155.13 | 1.2821D+04 |
| 0.75 | 735 | 258 | 3 / 271 | 278 | 278 | 271 | 818.61 | 1.6527D+04 |
| 1 | 1000 | 299 | 3 / 348 | 355 | 355 | 348 | 1063.91 | 2.1418D+04 |

Table 2: Performance of Algencon 3.0.0 without the acceleration process and SECO (Algorithms 2.1 and 3.1) with no backtracking applied to six instances of problem (27) with $n = 1000$, $m = 500$, $|I| \approx pn$, and $p \in \{0, 0.01, 0.1, 0.5, 0.75, 1\}$.

problem P_j in order to arrive to a point x^* with functional value $f(x^*) = f_{ij}$ and feasibility $\|h(x^*)\|_\infty = h_{ij}$. It is assumed that the metric t_{ij} is such that the smaller its value, the higher the performance of method M_i on problem P_j . Moreover, let t_j^{\min} denote the smallest among all the performance measurements required by each method that “found a solution” for problem P_j . In performance profiles, each method M_i is related to a curve

$$\Gamma_i(\tau) = \frac{\#\{j \in \{1, \dots, p\} \mid M_i \text{ found a solution for } P_j \text{ with } t_{ij} \leq \tau t_j^{\min}\}}{p},$$

where $\#\mathcal{S}$ denotes the cardinality of set \mathcal{S} . Let

$$f_j^{\min} = \min_{1 \leq i \leq q} \{f_{ij} \mid h_{ij} \leq \varepsilon^{\text{feas}}\}$$

and consider

$$\varepsilon_{ij} = \frac{f_{ij} - f_j^{\min}}{\max\{1, |f_j^{\min}|\}}. \quad (28)$$

For a given tolerance $\varepsilon_f > 0$, we say that *method M_i found a solution* to problem P_j if

$$h_{ij} \leq \varepsilon^{\text{feas}} \text{ and } \varepsilon_{ij} \leq \varepsilon_f. \quad (29)$$

In addition, we also say that method M_i found a solution to problem P_j if

$$h_{ij} \leq \varepsilon^{\text{feas}} \text{ and } f_{ij} \leq -f_\infty,$$

where f_∞ is a very large positive number. In this case, we assume the objective function is unbounded from below within the feasible region and any value of $f_{ij} \leq -f_\infty$ is considered a solution. In the numerical comparison, we considered the CPU time that a method M_i took on a problem P_j to find a point x^* that satisfies the method’s stopping criterion as the performance measurement t_{ij} . We arbitrarily set $\varepsilon_f = 10^{-4}$ and $f_\infty = 10^{10}$ (and $\varepsilon^{\text{feas}} = 10^{-8}$). This relatively loose value assigned to ε_f is essential to cancel the advantage in obtaining lower functional values of methods that are allowed to satisfy the bound constraints within a given tolerance (in contrast with Algencan that satisfies the bound constraints with no tolerance).

Figure 1 shows the results. In the figure, the x -axis correspond to $\log_{10}(\tau)$ while the y -axis corresponds to $\Gamma(\tau)$. For a given method, the values of $\Gamma(1)$ and $\Gamma(\infty)$, that are displayed in the figure, are commonly associated with the efficiency and the robustness of the method, respectively. Therefore, it is worth mentioning that the combination `PERTURBX0 = FALSE` and `NOBCKTRCKATALL = FALSE` is the more efficient and the most robust. It is also noticeable that the perturbation of the initial point (that apparently is “very small”) has a negative effect in the robustness of the method, independently of performing backtracking or not. Since the initial point is provided together with the problem definition, this appears to be a characteristic of the CUTEst collection’s considered problems; characteristic that was not observed in the problems considered in the previous subsections. The fact the most robust combinations being the ones that *do* perform backtracking is also in opposition with the performance observed in the problems of the previous subsections.

A final figure that should be mentioned is the number of times each combination satisfied at least one of the stopping criteria at Steps 1.1–1.2 of Algorithm 3.1 (within the imposed CPU time limit). Those numbers are, from top to down in the order the combinations are reported in Figure 1, 157, 142, 145, 131, respectively. Let us focus for a moment on the performance of Algorithm 3.1 with `PERTURBX0 = FALSE` and `NOBCKTRCKATALL = FALSE` that found an approximate stationary point in 157 problems (out of 190 problems and within the considered CPU time limit). This means that, ignoring the value of the objective function and considering for a moment that the goal of the method is to find a stationary point, the rate of success of the method is 83%. Considering that the set of problems being used has problems that are known to be infeasible, this is a relatively high rate of success, that is comparable to the one obtained by much more sophisticated methods recently introduced in the literature and that posses global and local convergence theories as well as rate of convergence results. The conclusion may be that, if such simple method as Algorithm 3.1 is able “to solve” 83% of the problems from the CUTEst collection with only equality constraints, better results should be expected

from more sophisticated methods.

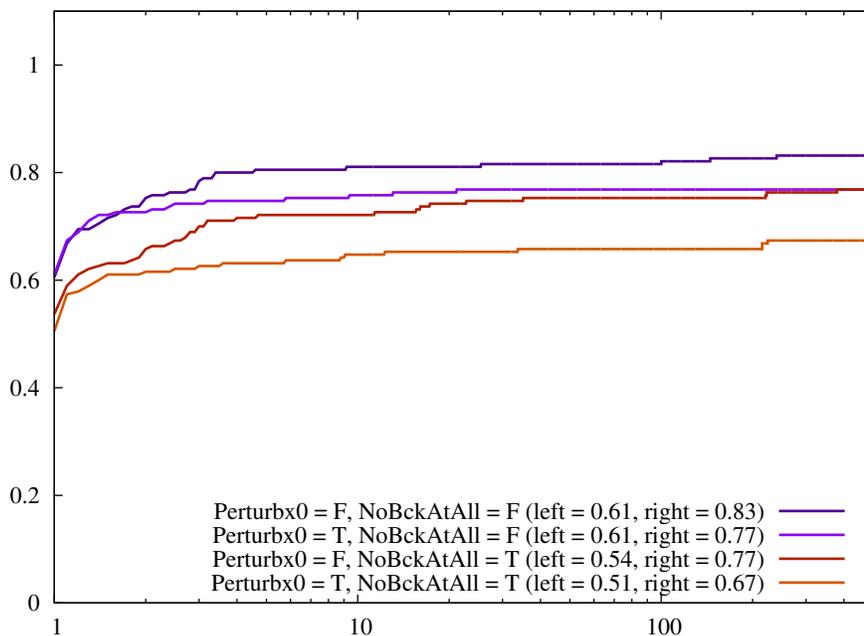


Figure 1: Comparison of Algorithm 3.1 with all four combinations of `PERTURBX0` and `NOBCKTRCKATALL` $\in \{\text{FALSE}, \text{TRUE}\}$.

Figure 2 shows a comparison (using performance profiles) between Algorithm 3.1 with `PERTURBX0 = FALSE` and `NOBCKTRCKATALL = FALSE` and Algencan. It is very clear that, for the considered set of problems and within the described comparison framework, the former method is more efficient and robust than Algencan. The comparison process that takes into account the objective function value at the final iterate says that both methods found different stationary points in 13 problems. Among these 13 problems, Algencan was declared as having failed in 8 problems; while Algorithm 3.1 was declared as having failed in the other 5 problems. This gives an advantage of $8 - 5 = 3$ problems or $3/190 \times 100\% \approx 1.5\%$ in the robustness index of Algorithm 3.1 that is 8 percentage points larger than the Algencan's robustness index; meaning that the superior robustness of Algorithm 3.1 does not depend on the comparison process but is due to the fact of Algorithm 3.1 having found stationary points more often than Algencan.

Numerical experiments in Sections 4.1 and 4.2 highlighted that perturbing the initial guess or avoiding the backtracking strategy may be profitable in some situations. However, numerical experiments with the CUTEst' test problems in Section 4.4.1 pointed out that the most efficient and robust version of Algorithm 3.1 is the one that does not perturb the initial guess and performs backtracking. This means that we have failed in concluding

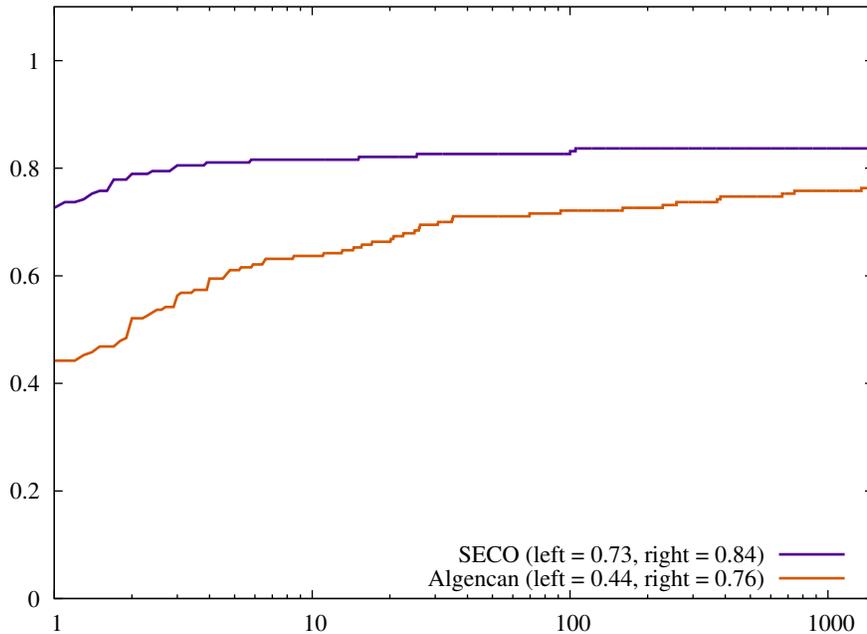


Figure 2: Comparison between Algorithm 3.1 with `PERTURBX0 = FALSE` and `NOBCK-TRCKATALL = FALSE` (identified as SECO in the graphic) and Algencan 3.0.0 considering *all* problems from CUTEst collection with only equality constraints.

whether these strategies should be considered or not in a first trial of solving a given problem; and that the interested user should play with these possibilities after a first unsuccessful trial of solving a problem at hand.

4.4.2 Problems with equality constraints and bound constraints

In this set of experiments we considered *all* the 283 problems from the CUTEst collection with at least one equality constraint, at least one bound constraint, and no inequality constraints (default dimensions and provided initial guesses were considered). Aiming to asset the reliability of the choice of lower- and upper-level constraints made in the development of the SECO algorithm, we compared SECO against Algencan (without considering the Algencan’s acceleration process; see [8] for details). In the latter case, equality constraints are penalized and bound-constrained subproblems are solved; while, in the former case, bound constraints are penalized and equality-constrained subproblems are solved. This is the qualitative difference between the two methods being compared and all other algorithmic parameters and implementation issues (like, for example, the used linear algebra subroutines) are the same.

As in the previous section, we used performance profiles to compare the two methods. In a first comparison, performance profiles were used as described in the previous section, i.e. considering the value of the objective function at the final iterate to determine whether a method solved a problem or not. Figure 3 shows the results. It is easy to see that Algencan appears as more efficient and robust than SECO. Two main factors determine this results: (a) the behaviour of SECO in a family of 45 problems and (b) a slightly larger tendency of Algencan to find solutions with lower functional value. We now analyze both factors in separate. The considered set of 283 problems from the CUTEst collection includes 45 problems (representing 16% of the problems) that are quadratic programming reformulations of linear complementarity problems (provided by Michael Ferris). In 11 out of this 45 problems, SECO presented a phenomenon named greediness in [12, 10] that may affect penalty and Lagrangian methods when the objective function takes very low values (perhaps going to $-\infty$) in the non-feasible region. In this case, iterates of the subproblems' solver may be attracted by undesired minimizers, especially at the first outer iterations, and overall convergence may fail to occur. The behaviour of SECO in those 11 problems reduces its robustness in $(11/283) \times 100\% \approx 4\%$. Independently of that, there are also 13 problems in which both methods found a feasible point but criterion (28–29) says that, due to the difference in the functional value, one of the methods succeeded while the other has failed. Among those 13 problems, in 9 cases it was declared that Algencan has found a solution and SECO has failed; while the opposite situation happened in the remaining 4 problems. The difference $9 - 4 = 5$ affects the robustness of SECO in approximately 2% and this fully explains the difference of 6% between the robustness indices of Algencan and SECO. If we eliminate those $45 + 13 = 58$ problems from the comparison, we obtain the performance profile depicted in Figure 4 (that is based in the remaining 225 problems only). This figure shows, in the considered subset of problems and having already mentioned that Algencan appears to show a slightly larger tendency to find lower functional values, that both methods are almost equally efficient and robust. In addition, note that the satisfaction of the KKT conditions (norm of the gradient of the Lagrangian and complementarity conditions) plays no role in this comparison. However, with respect to the satisfaction of the KKT conditions in the subset of 225 problems being considered, we can say that Algencan satisfied the KKT conditions in 150 problems while SECO satisfied the KKT conditions in 155 problems (always within the CPU time limit of 5 minutes). As a whole, this appears to be a surprising result considering the unorthodox choice of lower- and upper-level constraints made in the SECO algorithm.

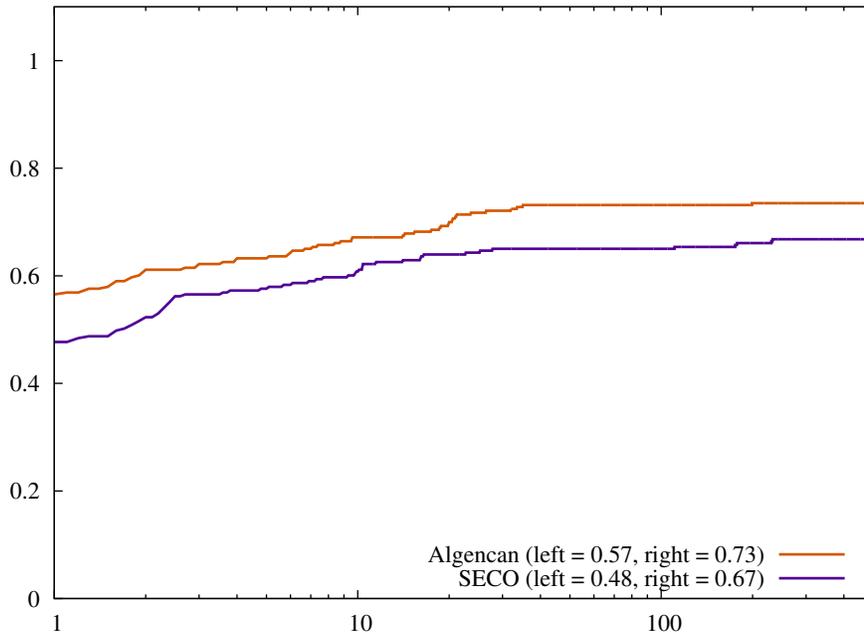


Figure 3: Comparison between SECO (Algorithms 2.1–3.1 with `PERTURBX0 = FALSE` and `NOBCKTRCKATALL = FALSE`) and Algencan 3.0.0 without acceleration, considering *all* problems from CUTEst collection with equality and bound constraints.

5 Conclusions

General Augmented Lagrangian methods [1, 9] are based on the partition of the constraints of the nonlinear programming problem into two sets; the first one is defined by

$$h(x) = 0 \text{ and } g(x) \leq 0 \quad (30)$$

and the second is given by

$$\underline{h}(x) = 0 \text{ and } \underline{g}(x) \leq 0, \quad (31)$$

where $h : \mathbb{R}^n \rightarrow \mathbb{R}^m$, $g : \mathbb{R}^n \rightarrow \mathbb{R}^p$, $\underline{h} : \mathbb{R}^n \rightarrow \mathbb{R}^m$, and $\underline{g} : \mathbb{R}^n \rightarrow \mathbb{R}^p$. The subproblems of the AL algorithm consist on the minimization of the Augmented Lagrangian defined by f , h , and g subject to the constraints (31).

The constraints (31) have been called “lower-level constraints”, “simple constraints”, or “subproblem constraints” in different papers. In most practical algorithms the constraints (31) are really simple, for example, when they define a box. However, in the present paper we show that the denomination “simple” for the constraints (31) may not be adequate because it may be attractive to consider that the constraints (31) are “more complicated” than the constraints (30).

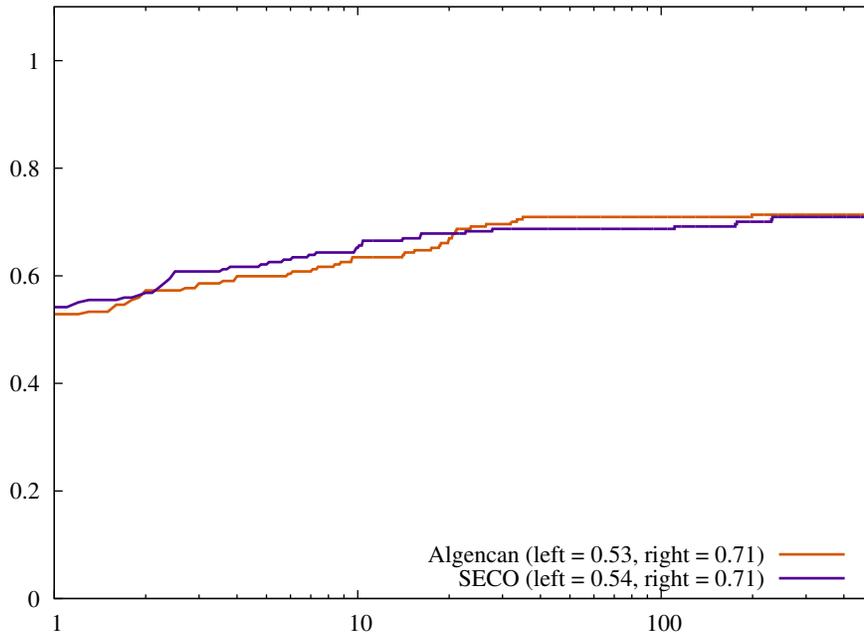


Figure 4: Comparison between SECO (Algorithms 2.1–3.1 with `PERTURBX0 = FALSE` and `NOBCKTRCKATALL = FALSE`) and Algencan 3.0.0 without acceleration, considering a *subset* of 225 problems from CUTEst collection with equality and bound constraints.

To fix ideas, suppose that the constraints (30) are defined by

$$h(x) = 0 \tag{32}$$

and the constraints (31) are

$$\ell \leq x \leq u. \tag{33}$$

Every constrained optimization problem can be expressed as the minimization of a function $f(x)$ subject to (32) and (33), using slack variables. The conventional wisdom is that, if we adopt the AL strategy, we should solve this problem by means of subproblems that minimize Augmented Lagrangians defined by f and h subject to (33). In this paper we explored the idea of solving the problem by means of successive minimizations of the Augmented Lagrangian defined by f and the functions $x - u$ and $\ell - x$, subject to the constraints (32). Of course we made our best to define a method for solving the equality-constrained subproblems by means of which the overall algorithm should be competitive. The resulting algorithm have been compared with the well-established software Algencan. The conclusions of the numerical study depend on the definition of “satisfactory solution”. If we adopt the criterion that a satisfactory solution is a point that satisfies approximately

the KKT conditions with high precision, we can establish that the new algorithm is more efficient than Algencan when the number of active bounds is small, but its relative superiority in terms of efficiency decreases as the number of active constraints at the solution increases. The considerations that explain this behavior are the following:

1. In the extreme case in which there are no active inequality constraints the algorithm essentially solves a smooth nonlinear system of equations using a safeguarded Newton-like method. It is not surprising that this type of procedure should be more efficient than the possibly painful process of increasing the penalty parameter several times in order to minimize Augmented Lagrangians defined by f and h .
2. If many inequality constraints are active at the solution, even the algorithm for minimizing with equality constraints being efficient, several outer iterations are generally needed, because many bound constraints are not satisfied after each outer iteration. Moreover, the nonlinear system that represents the Lagrange conditions for that subproblem is semismooth, but not smooth, decreasing partially the efficiency of Newton-like algorithms. On the other hand, the classical implementation of Algencan deals efficiently with the active bound constraints, avoiding useless evaluations of f and h outside the feasible box.

Nevertheless, if we adopt the criterion that algorithms should be compared also with respect to the objective function value achieved at the end of execution, Algencan seems to be better than SECO. This can be explained by the opportunistic line-search strategies used by the bound-constrained solver employed by Algencan, which actively searches low values of the objective function with procedures that are not necessarily linked to worst-case convergence analysis.

We would like to finish this paper with strong statements concerning the relative efficiency or robustness of the new method with respect to existing ones. However, we are convinced that such statements are not possible in nonlinear optimization problems. There exist an infinitely large number of geometrical combinations, structures, and dimensions that make it impossible to claim universal superiority. Therefore, we conclude with the cautious claim that the SECO approach may be a useful tool when the problems are formulated with the constraints “ $h(x) = 0$ plus bounds” and the number of expected active bounds at the solution is moderate.

References

- [1] R. Andreani, E. G. Birgin, J. M. Martínez, and M. L. Schuverdt, On Augmented Lagrangian methods with general lower-level constraints, *SIAM Journal on Optimization* 18, pp. 1286–1309, 2008.

- [2] R. Andreani, G. Haeser, and J. M. Martínez, On sequential optimality conditions for smooth constrained optimization, *Optimization* 60, pp. 627–641, 2011.
- [3] R. Andreani, G. Haeser, M. L. Schuverdt, and P. J. S. Silva, A relaxed constant positive linear dependence constraint qualification and applications, *Mathematical Programming* 135, pp. 255–273, 2012.
- [4] R. Andreani, G. Haeser, M. L. Schuverdt, and P. J. S. Silva, Two new weak constraint qualifications and applications, *SIAM Journal on Optimization* 22, pp. 1109–1135, 2012.
- [5] R. Andreani, J. M. Martínez, A. Ramos, and P. J. S. Silva, A cone-continuity constraint qualification and algorithmic consequences, *SIAM Journal on Optimization* 26, pp. 96–110, 2016.
- [6] R. Andreani, J. M. Martínez, and M. L. Schuverdt, On the relation between Constant Positive Linear Dependence Condition and Quasinormality Constraint Qualification, *Journal of Optimization Theory and Applications* 125, pp. 473–485, 2005.
- [7] N. Banihashemi and C. Y. Kaya, Inexact Restoration for Euler Discretization of Box-Constrained Optimal Control Problems, *Journal of Optimization Theory and Applications* 156, pp. 726–760, 2013.
- [8] E. G. Birgin and J. M. Martínez, Improving ultimate convergence of an Augmented Lagrangian method, *Optimization Methods and Software* 23, pp. 177–195, 2008.
- [9] E. G. Birgin and J. M. Martínez, *Practical Augmented Lagrangian Methods for Constrained Optimization*, SIAM, Philadelphia, 2014.
- [10] E. G. Birgin, E. V. Castalani, A. L. M. Martinez, and J. M. Martínez, Outer Trust-Region method for Constrained Optimization, *Journal of Optimization Theory and Applications* 150, pp. 142–155, 2011.
- [11] L. F. Bueno, G. Haeser, and J. M. Martínez, A Flexible Inexact Restoration Method for Constrained Optimization, *Journal of Optimization Theory and Applications* 165, pp. 188–208, 2015.
- [12] E. V. Castalani, A. L. Martinez, J. M. Martínez, and B. F. Svaiter, Addressing the greediness phenomenon in Nonlinear Programming by means of Proximal Augmented Lagrangians, *Computational Optimization and Applications* 46, 229–245, 2010.
- [13] A. R. Conn, N. I. M. Gould, and Ph. L. Toint, *Lancelot: A Fortran package for large scale nonlinear optimization*, Springer-Verlag, Berlin, 1992.

- [14] E. D. Dolan and J. J. Moré, Benchmarking optimization software with performance profiles, *Mathematical Programming* 91, pp. 201–213, 2002.
- [15] A. Fischer and A. Friedlander, A new line search inexact restoration approach for nonlinear programming, *Computational Optimization and Applications* 46, pp. 333–346, 2010.
- [16] R. Fletcher, *Practical Methods of Optimization*, 2nd edition, Wiley, 2000.
- [17] J. B. Francisco, J. M. Martínez, L. Martínez, and F. Pisnitchenko, Inexact Restoration method for minimization problems arising in electronic structure calculations, *Computational Optimization and Applications* 50, pp. 555–590, 2011.
- [18] G. H. Golub and C. F. Van Loan, *Matrix Computations* (third edition), The Johns Hopkins University Press, Baltimore and London, 1996.
- [19] C. C. Gonzaga, E. W. Karas, and M. Vanti, A globally convergent filter method for Nonlinear Programming, *SIAM Journal on Optimization* 14, pp. 646–669, 2004.
- [20] N. I. M. Gould, D. Orban, and Ph. L. Toint, CUTEst: a constrained and unconstrained testing environment with safe threads for mathematical optimization, *Computational Optimization and Applications* 60, pp. 545–557, 2014.
- [21] E. W. Karas, C. C. Gonzaga, and A. A. Ribeiro, Local convergence of filter methods for equality constrained non-linear programming, *Optimization* 59, pp. 1153–1171, 2010.
- [22] E. W. Karas, E. A. Pilotta, and A. A. Ribeiro, Numerical comparison of merit function with filter criterion in inexact restoration algorithms using Hard-Spheres Problems, *Computational Optimization and Applications* 44, pp. 427–441, 2009.
- [23] C. Y. Kaya, Inexact Restoration for Runge-Kutta discretization of Optimal Control problems, *SIAM Journal on Numerical Analysis* 48, pp. 1492–1517, 2010.
- [24] C. Y. Kaya and J. M. Martínez, Euler discretization and Inexact Restoration for Optimal Control, *Journal of Optimization Theory and Application* 134, pp. 191–206, 2007.
- [25] O. L. Mangasarian and S. Fromovitz, The Fritz-John necessary optimality conditions in presence of equality and inequality constraints, *Journal of Mathematical Analysis and Applications* 17 pp. 37–47, 1967.
- [26] J. M. Martínez, Inexact-Restoration method with Lagrangian tangent decrease and new merit function for nonlinear programming, *Journal of Optimization Theory and Application* 111, 39–58, 2001.

- [27] J. M. Martínez and E. A. Pilotta, Inexact-Restoration algorithms for constrained optimization, *Journal of Optimization Theory and Application* 104, pp. 135–163, 2000.
- [28] B. A. Murtagh and M. A. Saunders, Large-scale linearly constrained optimization, *Mathematical Programming* 14, pp. 41–72, 1978.
- [29] J. Nocedal and S. J. Wright, *Numerical Optimization*, 2nd edition, Springer, 2006.
- [30] L. Qi and J. Sun, A nonsmooth version of Newton’s method, *Mathematical Programming* 58, pp. 353–367, 1993.
- [31] L. Qi and Z. Wei, On the constant linear dependence condition and its application to SQP methods, *SIAM Journal on Optimization* 10, pp. 963–981, 2000.
- [32] A. Wächter and L. T. Biegler, On the implementation of an interior-point filter line-search algorithm for large-scale nonlinear programming, *Mathematical Programming* 106, pp. 25-57, 2006.
- [33] HSL(2013), *A collection of Fortran codes for large scale scientific computation*, <http://www.hsl.rl.ac.uk>.