

An adaptive cubic regularization algorithm for nonconvex optimization with convex constraints and its function-evaluation complexity

C. CARTIS

School of Mathematics, University of Edinburgh, Edinburgh EH9 3JZ, United Kingdom
coralia.cartis@ed.ac.uk

N. I. M. GOULD*

Computational Science and Engineering Department, Rutherford Appleton Laboratory,
Chilton OX11 0QX, United Kingdom

*Corresponding author: nick.gould@stfc.ac.uk

AND

PH. L. TOINT

Department of Mathematics, FUNDP-University of Namur, 61, Rue de Bruxelles,
B-5000 Namur, Belgium

philippe.toint@fundp.ac.be

[Received on 24 April 2010; revised on 14 July 2011]

The adaptive cubic regularization algorithm described in Cartis *et al.* (2009, Adaptive cubic regularization methods for unconstrained optimization. Part I: motivation, convergence and numerical results. *Math. Program.*, **127**, 245–295; 2010, Adaptive cubic regularisation methods for unconstrained optimization. Part II: worst-case function- and derivative-evaluation complexity [online]. *Math. Program.*, DOI: 10.1007/s10107-009-0337-y) is adapted to the problem of minimizing a nonlinear, possibly non-convex, smooth objective function over a convex domain. Convergence to first-order critical points is shown under standard assumptions, without any Lipschitz continuity requirement on the objective's Hessian. A worst-case complexity analysis in terms of evaluations of the problem's function and derivatives is also presented for the Lipschitz continuous case and for a variant of the resulting algorithm. This analysis extends the best-known bound for general unconstrained problems to nonlinear problems with convex constraints.

Keywords: nonlinear optimization; convex constraints; cubic regularization/regularization; numerical algorithms; global convergence; worst-case complexity.

1. Introduction

Adaptive cubic regularization has recently returned to the forefront of smooth nonlinear optimization as a possible alternative to more standard globalization techniques for unconstrained optimization. Methods of this type—initiated independently by Griewank (1981), Nesterov & Polyak (2006) and Weiser *et al.* (2007)—are based on the observation that a second-order model involving a cubic term can be constructed that overestimates the objective function when the latter has a Lipschitz continuous Hessian and a model parameter is chosen large enough. In Cartis *et al.* (2011a), we have proposed updating the parameter so that it merely estimates a local Lipschitz constant of the Hessian, as well as using

approximate model Hessians and approximate model minimizers, which makes this suitable for large-scale problems. These adaptive regularization methods are not only globally convergent to first- and second-order critical points with fast asymptotic speed (Nesterov & Polyak, 2006; Cartis *et al.*, 2011a), but also—unprecedentedly—enjoy better worst-case global complexity bounds than steepest descent methods (Nesterov & Polyak, 2006; Cartis *et al.*, 2011b), Newton’s method and trust-region methods (Cartis *et al.*, 2010). Furthermore, preliminary numerical experiments with basic implementations of these techniques and of the trust region show encouraging performance of the cubic regularization approach (Cartis *et al.*, 2011a).

Extending the approach to more general optimization problems is therefore attractive, as one may hope that some of the qualities of the unconstrained methods can be transferred to a broader framework. Nesterov (2006) has considered the extension of his cubic regularization method to problems with a smooth convex objective function and convex constraints. In this paper we consider the extension of the adaptive cubic regularization methods to the case where minimization is subject to convex constraints, but the smooth objective function is no longer assumed to be convex. The new algorithm is strongly inspired by the unconstrained adaptive cubic regularization methods (Cartis *et al.*, 2011a,b) and by the trust-region projection methods for the same constrained problem class that are fully described in of Conn *et al.* (2000, Chapter 12). In particular, it makes significant use of the specialized first-order criticality measure developed by Conn *et al.* (1993) for the latter context. Firstly, global convergence to first-order critical points is shown under mild assumptions on the problem class for a generic adaptive cubic regularization framework that only requires a Cauchy-like decrease in the (constrained) model subproblem. The latter can be efficiently computed using a generalized Goldstein linesearch, suitable for the cubic model, provided projections onto the feasible set are inexpensive to calculate. The associated worst-case global complexity—or equivalently, the total number of objective function and gradient evaluations—required by this generic cubic regularization approach to reach approximate first-order optimality matches, in order, that of steepest descent for unconstrained (nonconvex) optimization.

However, in order to improve the local and global rate of convergence of the algorithm, it is necessary to advance beyond the Cauchy point when minimizing the model. To this end we propose an adaptive cubic regularization variant that under certain assumptions on the algorithm, can be proved to satisfy the desirable global evaluation complexity bound of its unconstrained counterpart, which, as mentioned in the first paragraph, is better than for steepest descent methods. As in the unconstrained case we do not rely on global model minimization and are content with only sequential line minimizations of the model provided they ensure descent at each (inner) step. Possible descent paths of this type are suggested, though more work is needed to transform these ideas into a computationally efficient model solution procedure. Solving the (constrained) subproblem relies on the assumption that these piecewise linear paths are uniformly bounded, which still requires both practical and theoretical validation.

Our complexity analysis here, in terms of the function-evaluations count, does not cover the total computational cost of solving the problem as it ignores the cost of solving the (constrained) subproblem. Note, however, that though the latter may be NP-hard computationally (Vavasis, 1991), it does not require any additional function evaluations. Furthermore, for many examples, the cost of these (black-box) evaluations significantly dominates that of the internal computations performed by the algorithm. Even so, effective step calculation is crucial for the practical computational efficiency of the algorithm and will be given priority consideration in our future work.

The paper is organized as follows. Section 2 describes the constrained problem more formally as well as the new adaptive regularization algorithm for it, while Section 3 presents the associated convergence theory (to first-order critical points). We then discuss a worst-case function-evaluation complexity result

for the new algorithm and an improved result for a cubic regularization variant in Section 4. Finally some conclusions are presented in Section 5.

2. The new algorithm

We consider the numerical solution of the constrained nonlinear optimization problem

$$\min_{x \in \mathcal{F}} f(x), \quad (2.1)$$

where we assume that $f : \mathfrak{R}^n \rightarrow \mathfrak{R}$ is twice continuously differentiable, possibly nonconvex and bounded below on the closed, convex and nonempty feasible domain $\mathcal{F} \subseteq \mathfrak{R}^n$.

Our algorithm for solving this problem follows the broad lines of the projection-based trust-region algorithm of in [Conn *et al.* \(2000, Chapter 12\)](#) with adaptations necessary to replace the trust-region globalization mechanism by a cubic regularization of the type analysed in [Cartis *et al.* \(2011a\)](#). At an iterate x_k within the feasible region \mathcal{F} , a cubic model of the form

$$m_k(x_k + s) = f(x_k) + \langle g_k, s \rangle + \frac{1}{2} \langle s, B_k s \rangle + \frac{1}{3} \sigma_k \|s\|^3 \quad (2.2)$$

is defined, where $\langle \cdot, \cdot \rangle$ denotes the Euclidean inner product, $g_k \stackrel{\text{def}}{=} \nabla_x f(x_k)$, B_k is a symmetric matrix hopefully approximating the objective's Hessian $H(x_k) \stackrel{\text{def}}{=} \nabla_{xx} f(x_k)$, σ_k is a positive regularization parameter and $\|\cdot\|$ stands for the Euclidean norm. The step s_k from x_k is then defined in two stages. The first stage is to compute a *generalized Cauchy point* x_k^{GC} such that x_k^{GC} approximately minimizes the model (2.2) along the Cauchy arc defined by the projection onto \mathcal{F} of the negative gradient path; that is

$$\{x \in \mathcal{F} \mid x = P_{\mathcal{F}}[x_k - t g_k], t \geq 0\},$$

where we define $P_{\mathcal{F}}$ to be the (unique) orthogonal projector onto \mathcal{F} . The approximate minimization is carried out using a generalized Goldstein-like linesearch on the arc, as explained in [Conn *et al.* \(2000, Section 12.1\)](#). In particular, $x_k^{\text{GC}} = x_k + s_k^{\text{GC}}$ is determined such that

$$x_k^{\text{GC}} = P_{\mathcal{F}} \left[x_k - t_k^{\text{GC}} g_k \right] \quad \text{for some } t_k^{\text{GC}} > 0, \quad (2.3)$$

and

$$m_k \left(x_k^{\text{GC}} \right) \leq f(x_k) + \kappa_{\text{ubs}} \langle g_k, s_k^{\text{GC}} \rangle \quad (2.4)$$

and either

$$m_k \left(x_k^{\text{GC}} \right) \geq f(x_k) + \kappa_{\text{lbs}} \langle g_k, s_k^{\text{GC}} \rangle \quad (2.5)$$

or

$$\left\| P_{T(x_k^{\text{GC}})}[-g_k] \right\| \leq \kappa_{\text{epp}} |\langle g_k, s_k^{\text{GC}} \rangle|, \quad (2.6)$$

where the three constants satisfy

$$0 < \kappa_{\text{ubs}} < \kappa_{\text{lbs}} < 1 \quad \text{and} \quad \kappa_{\text{epp}} \in \left(0, \frac{1}{2} \right), \quad (2.7)$$

and where $T(x)$ is the tangent cone to \mathcal{F} at x . The conditions (2.4) and (2.5) are the familiar Goldstein linesearch conditions adapted to our search along the Cauchy arc, while (2.6) is there to handle the case where this arc ends before condition (2.5) is ever satisfied. Once the generalized Cauchy point x_k^{GC} is computed (which can be done by a suitable search on $t_k^{\text{GC}} > 0$ inspired by Conn *et al.*, 2000, Algorithm 12.2.2, and discussed below), any step s_k such that

$$x_k^+ \stackrel{\text{def}}{=} x_k + s_k \in \mathcal{F}$$

and such that the model value at x_k^+ is below that obtained at x_k^{GC} , is acceptable.

Given the step s_k , the trial point x_k^+ is known and the value of the objective function at this point computed. If the ratio

$$\rho_k = \frac{f(x_k) - f(x_k^+)}{f(x_k) - m_k(x_k^+)} \quad (2.8)$$

of the achieved reduction in the objective function compared to the predicted model reduction is larger than some constant $\eta_1 > 0$, then the trial point is accepted as the next iterate and the regularization parameter σ_k is essentially unchanged or decreased, while the trial point is rejected and σ_k increased if $\rho_k < \eta_1$. Fortunately, the undesirable situation where the trial point is rejected cannot persist since σ_k eventually becomes larger than some local Lipschitz constant associated with the Hessian of the objective function (assuming it exists), which in turn guarantees that $\rho_k \geq 1$, as shown in Griewank (1981), Nesterov & Polyak (2006) or Cartis *et al.* (2011a).

We now state our Adaptive Regularization using Cubics for CONvex Constraints (COCARC).

Algorithm 2.1. Adaptive Regularization with Cubics for Convex Constraints (COCARC).

Step 0. Initialization. An initial point $x_0 \in \mathcal{F}$ and an initial regularization parameter $\sigma_0 > 0$ are given. Compute $f(x_0)$ and set $k = 0$.

Step 1. Determination of the generalized Cauchy point. If x_k is first-order critical, terminate the algorithm. Otherwise perform the following iteration.

Step 1.0. Initialization. Define the model (2.2), choose $t_0 > 0$ and set

$$t_{\min} = 0, t_{\max} = \infty \text{ and } j = 0.$$

Step 1.1. Compute a point on the projected-gradient path. Set $x_{k,j} = P_{\mathcal{F}}[x_k - t_j g_k]$ and evaluate $m_k(x_{k,j})$.

Step 1.2. Check for the stopping conditions. If (2.4) is violated then set $t_{\max} = t_j$ and go to Step 1.3. Otherwise, if (2.5) and (2.6) are violated, set $t_{\min} = t_j$ and go to Step 1.3. Otherwise set $x_k^{\text{GC}} = x_{k,j}$ and go to Step 2.

Step 1.3. Find a new value of the arc parameter. If $t_{\max} = \infty$, set $t_{j+1} = 2t_j$. Otherwise set $t_{j+1} = \frac{1}{2}(t_{\min} + t_{\max})$. Increment j by 1 and go to Step 1.2.

Step 2. Step calculation. Compute a step s_k and a trial point $x_k^+ \stackrel{\text{def}}{=} x_k + s_k \in \mathcal{F}$ such that

$$m_k(x_k^+) \leq m_k(x_k^{\text{GC}}). \quad (2.9)$$

Step 3. Acceptance of the trial point. Compute $f(x_k^+)$ and the ratio (2.8). If $\rho_k \geq \eta_1$ then define $x_{k+1} = x_k + s_k$; otherwise define $x_{k+1} = x_k$.

Step 4. Regularization parameter update. Set

$$\sigma_{k+1} \in \begin{cases} (0, \sigma_k] & \text{if } \rho_k \geq \eta_2, \\ [\sigma_k, \gamma_1 \sigma_k] & \text{if } \rho_k \in [\eta_1, \eta_2), \\ [\gamma_1 \sigma_k, \gamma_2 \sigma_k] & \text{if } \rho_k < \eta_1. \end{cases}$$

Increment k by 1 and go to Step 1.

As in [Cartis *et al.* \(2011a\)](#) the constants η_1 , η_2 , γ_1 and γ_2 are given and satisfy the conditions

$$0 < \eta_1 \leq \eta_2 < 1 \quad \text{and} \quad 1 < \gamma_1 \leq \gamma_2. \quad (2.10)$$

As for trust-region algorithms we say that iteration k is successful whenever $\rho_k \geq \eta_1$ (and thus $x_{k+1} = x_k^+$) and very successful whenever $\rho_k \geq \eta_2$, in which case, additionally, $\sigma_{k+1} \leq \sigma_k$. We denote the index set of all successful and very successful iterations by \mathcal{S} .

As mentioned above, our technique for computing the generalized Cauchy point is inspired by the Goldstein linesearch scheme, but it is most likely that techniques based on Armijo-like backtracking (see [Sartenaer, 1993](#)) or on successive exploration of the active faces of \mathcal{F} along the Cauchy arc (see [Conn *et al.*, 1988](#)) are also possible, the latter being practical when \mathcal{F} is a polyhedron.

3. Global convergence to first-order critical points

We now consider the global convergence properties of Algorithm COCARC and show in this section that all the limit points of the sequence of its iterates must be first-order critical points for problem (2.1). Our analysis will be based on the first-order criticality measure at $x \in \mathcal{F}$ given by

$$\chi(x) \stackrel{\text{def}}{=} \left| \min_{x+d \in \mathcal{F}, \|d\| \leq 1} \langle \nabla_x f(x), d \rangle \right|, \quad (3.1)$$

(see [Conn *et al.*, 1993](#)) and define $\chi_k \stackrel{\text{def}}{=} \chi(x_k)$. We say that x_* is a first-order critical point for (2.1) if $\chi(x_*) = 0$ (see [Conn *et al.*, 2000](#) Theorem 12.1.6).

For our analysis, we consider the following assumptions.

AS1. The feasible set \mathcal{F} is closed, convex and nonempty.

AS2. The function f is twice continuously differentiable on the (open and convex) set $\hat{\mathcal{F}}_0 = \{x : \|x - y\| < \delta \text{ for some } y \in \mathcal{F}_0\}$ for given $\delta \in (0, 1)$ and where $\mathcal{F}_0 \subseteq \mathcal{F}$ is the closed, convex hull of x_0 and the iterates $x_k + s_k$, $k \geq 0$.

AS3a. The function f is bounded below by f_{low} on \mathcal{F}_0 .

AS3b. The set \mathcal{F}_0 is bounded.

AS4. There exist constants $\kappa_H > 1$ and $\kappa_B > 1$ such that

$$\|H(x)\| \leq \kappa_H \quad \text{for all } x \in \mathcal{F}_0, \quad \text{and} \quad \|B_k\| \leq \kappa_B \quad \text{for all } k \geq 0. \quad (3.2)$$

Note that AS3b and AS2 imply AS3a, but some results will only require the weaker condition AS3a.

Suppose that AS1 and AS2 hold, and let $x \in \mathcal{F}_0$. For $t > 0$, let

$$x(t) \stackrel{\text{def}}{=} P_{\mathcal{F}}[x - t \nabla_x f(x)] \quad \text{and} \quad \theta(x, t) \stackrel{\text{def}}{=} \|x(t) - x\|, \quad (3.3)$$

while, for $\theta > 0$,

$$\chi(x, \theta) \stackrel{\text{def}}{=} \left| \min_{x+d \in \mathcal{F}, \|d\| \leq \theta} \langle \nabla_x f(x), d \rangle \right|, \quad (3.4)$$

and

$$\pi(x, \theta) \stackrel{\text{def}}{=} \frac{\chi(x, \theta)}{\theta}. \quad (3.5)$$

Some already-known properties of the projected gradient path and the above variants of the criticality measure (3.1) are given next and will prove useful in what follows.

LEMMA 3.1

1. [Conn *et al.*, 2000] Suppose that AS1 and AS2 hold and let $x \in \mathcal{F}_0$ and $t > 0$ such that $\theta > 0$. Then

- (i) [Theorem 3.2.8] $\theta(x, t)$, $\chi(x, \theta)$ and $\pi(x, \theta)$ are continuous with respect to their two arguments,
- (ii) [Theorem 12.1.3] $\theta(x, t)$ is nondecreasing with respect to t ,
- (iii) [Theorem 12.1.4] the point $x(t) - x$ is a solution of problem

$$\min_{x+d \in \mathcal{F}, \|d\| \leq \theta} \langle \nabla_x f(x), d \rangle, \quad (3.6)$$

where $\theta = \|x(t) - x\|$,

- (iv) [Theorem 12.1.5(i), (ii)] $\chi(x, \theta)$ is nondecreasing and $\pi(x, \theta)$ is nonincreasing with respect to θ ,
- (v) [Theorem 12.1.5 (iii)] for any d such that $x + d \in \mathcal{F}$, the inequality

$$\chi(x, \theta) \leq |\langle \nabla_x f(x), d \rangle| + 2\theta \|P_{T(x+d)}[-\nabla_x f(x)]\| \quad (3.7)$$

holds for all $\theta > \|d\|$.

2. [Hiriart-Urruty and Lemařechal, 1993, Proposition 5.3.5] For any $x \in \mathcal{F}$ and $d \in \mathfrak{R}^n$, the following limit holds:

$$\lim_{\alpha \rightarrow 0^+} \frac{P_{\mathcal{F}}(x + \alpha d) - x}{\alpha} = P_{T(x)}[d]. \quad (3.8)$$

The following result is a consequence of the above properties of the criticality measure (3.1) and its variants.

LEMMA 3.2 Suppose that AS1 and AS2 hold. For $x \in \mathcal{F}_0$, $t > 0$ and $\theta > 0$, recall the measures (3.3), (3.4) and (3.5), and let

$$\pi_k^{\text{GC}} \stackrel{\text{def}}{=} \pi \left(x_k, \|s_k^{\text{GC}}\| \right) \quad \text{and} \quad \pi_k^+ \stackrel{\text{def}}{=} \pi(x_k, \|s_k\|), \quad (3.9)$$

where $s_k^{\text{GC}} \stackrel{\text{def}}{=} x_k^{\text{GC}} - x_k$. If $\|s_k^{\text{GC}}\| \geq 1$ then

$$\chi \left(x_k, \|s_k^{\text{GC}}\| \right) \geq \chi_k \geq \pi_k^{\text{GC}}, \quad (3.10)$$

while if $\|s_k^{\text{GC}}\| \leq 1$ then

$$\pi_k^{\text{GC}} \geq \chi_k \geq \chi \left(x_k, \|s_k^{\text{GC}}\| \right). \quad (3.11)$$

Similarly, if $\|s_k\| \geq 1$ then

$$\chi(x_k, \|s_k\|) \geq \chi_k \geq \pi_k^+, \quad (3.12)$$

while if $\|s_k\| \leq 1$ then

$$\pi_k^+ \geq \chi_k \geq \chi(x_k, \|s_k\|). \quad (3.13)$$

Moreover,

$$-\langle g_k, s_k^{\text{GC}} \rangle = \chi \left(x_k, \|s_k^{\text{GC}}\| \right) \geq 0, \quad (3.14)$$

$$\chi_k \leq \chi \left(x_k, \|s_k^{\text{GC}}\| \right) + 2 \left\| P_{T(x_k^{\text{GC}})}[-g_k] \right\| \quad (3.15)$$

and

$$\theta(x, t) \geq t \left\| P_{T(x(t))}[-\nabla_x f(x)] \right\|. \quad (3.16)$$

Proof. The inequalities (3.10) and (3.11) follow from the identity

$$\chi_k = \chi(x_k, 1), \quad (3.17)$$

(3.5) and Lemma 3.1(iv). Precisely the same arguments give (3.12) and (3.13) as well since the definition of s_k^{GC} was not used in the above inequalities. To show (3.14), apply Lemma 3.1(iii) with $t = t_k^{\text{GC}}$, which gives $\theta = \|s_k^{\text{GC}}\|$, and recalling the definition of (3.4), also

$$\left| \langle g_k, s_k^{\text{GC}} \rangle \right| = \chi \left(x_k, \|s_k^{\text{GC}}\| \right). \quad (3.18)$$

It remains to show that $|\langle g_k, s_k^{\text{GC}} \rangle| = -\langle g_k, s_k^{\text{GC}} \rangle$, which follows from the monotonicity of the projection operator, namely, we have

$$\left\langle x_k - t_k^{\text{GC}} g_k - x \left(t_k^{\text{GC}} \right), x_k - x \left(t_k^{\text{GC}} \right) \right\rangle \leq 0,$$

or equivalently,

$$\left\langle g_k, s_k^{\text{GC}} \right\rangle \leq -\frac{1}{t_k^{\text{GC}}} \left\| x_k - x(t_k^{\text{GC}}) \right\|^2 \leq 0.$$

Next, (3.15) results from (3.10) if $\|s_k^{\text{GC}}\| \geq 1$; else, when $\|s_k^{\text{GC}}\| < 1$, (3.15) follows by letting $x = x_k$, $\theta = 1$ and $d = s_k^{\text{GC}}$ in (3.7) and employing (3.18). We are left with proving (3.16). We first note that if $u(x, t) = x(t) - x$ then $\theta(x, t) = \|u(x, t)\|$ and, denoting the right directional derivative by d/dt_+ , we see that

$$\frac{d\theta}{dt_+}(x, t) = \frac{\left\langle \frac{du(x, t)}{dt_+}, u(x, t) \right\rangle}{\|u(x, t)\|} = \frac{\langle P_{T(x(t))}[-\nabla_x f(x)], u(x, t) \rangle}{\theta(x, t)}, \quad (3.19)$$

where to deduce the second equality, we used (3.8) with $x = x(t)$ and $d = -\nabla_x f(x)$. Moreover,

$$u(x, t) = -t\nabla_x f(x) - [x - t\nabla_x f(x) - x(t)] \stackrel{\text{def}}{=} -t\nabla_x f(x) - z(x, t) \quad (3.20)$$

and because of the definition of $x(t)$, $z(x, t)$ must belong to $N(x(t))$, the normal cone to \mathcal{F} at $x(t)$, which by definition, comprises all directions w such that $\langle w, y - x(t) \rangle \leq 0$ for all $y \in \mathcal{F}$. Thus, since this cone is the polar of $T(x(t))$, we deduce that

$$\langle P_{T(x(t))}[-\nabla_x f(x)], z(x, t) \rangle \leq 0. \quad (3.21)$$

We now obtain, successively using (3.19), (3.20) and (3.21), that

$$\begin{aligned} \theta(x, t) \frac{d\theta}{dt_+}(x, t) &= \langle P_{T(x(t))}[-\nabla_x f(x)], u(x, t) \rangle \\ &= \langle P_{T(x(t))}[-\nabla_x f(x)], -t\nabla_x f(x) - z(x, t) \rangle \\ &= t \langle -\nabla_x f(x), P_{T(x(t))}[-\nabla_x f(x)] \rangle - \langle P_{T(x(t))}[-\nabla_x f(x)], z(x, t) \rangle \\ &\geq t \|P_{T(x(t))}[-\nabla_x f(x)]\|^2. \end{aligned} \quad (3.22)$$

But (3.19) and the Cauchy–Schwarz inequality also imply that

$$\frac{d\theta}{dt_+}(x, t) \leq \|P_{T(x(t))}[-\nabla_x f(x)]\|.$$

Combining this last bound with (3.22) finally yields (3.16) as desired. \square

We complete our analysis of the criticality measures by considering the Lipschitz continuity of the measure $\chi(x)$. We start by proving the following lemma. This result extends [Mangasarian & Rosen \(1964, Lemma 1\)](#) by allowing a general, possibly implicit, expression of the feasible set.

LEMMA 3.3 Suppose that AS1 holds and define

$$\phi(x) \stackrel{\text{def}}{=} \min_{x+d \in \mathcal{F}, \|d\| \leq 1} \langle g, d \rangle$$

for $x \in \mathfrak{X}^n$ and some vector $g \in \mathfrak{X}^n$. Then $\phi(x)$ is a proper convex function on

$$\mathcal{F}_1 \stackrel{\text{def}}{=} \{x \in \mathfrak{X}^n \mid (\mathcal{F} - x) \cap \overline{\mathcal{B}} \neq \emptyset\}, \quad (3.23)$$

where $\overline{\mathcal{B}}$ is the closed Euclidean unit ball.

Proof. The result is trivial if $g = 0$. Assume, therefore, that $g \neq 0$. We first note that the definition of \mathcal{F}_1 ensures that the feasible set of $\phi(x)$ is nonempty and, therefore, that the parametric minimization problem defining $\phi(x)$ is well defined for any $x \in \mathcal{F}_1$. Moreover, the minimum is always attained because of the constraint $\|d\| \leq 1$, and so $-\infty < \phi(x)$ for all $x \in \mathcal{F}_1$. Hence, $\phi(x)$ is proper in \mathcal{F}_1 . To show that $\phi(x)$ is convex on (the convex set) \mathcal{F}_1 , let $x_1, x_2 \in \mathcal{F}_1$, and let $d_1, d_2 \in \mathfrak{R}^n$ be such that

$$\phi(x_1) = \langle g, d_1 \rangle \quad \text{and} \quad \phi(x_2) = \langle g, d_2 \rangle.$$

Also let $\lambda \in [0, 1]$, $x_0 \stackrel{\text{def}}{=} \lambda x_1 + (1 - \lambda)x_2$ and $d_0 \stackrel{\text{def}}{=} \lambda d_1 + (1 - \lambda)d_2$. Let us show that d_0 is feasible for the $\phi(x_0)$ problem. Since d_1 and d_2 are feasible for the $\phi(x_1)$ and $\phi(x_2)$ problems, respectively, and since $\lambda \in [0, 1]$, we have that $\|d_0\| \leq 1$. To show $x_0 + d_0 \in \mathcal{F}$, we have

$$x_0 + d_0 = \lambda(x_1 + d_1) + (1 - \lambda)(x_2 + d_2) \in \lambda\mathcal{F} + (1 - \lambda)\mathcal{F} \subseteq \mathcal{F},$$

where we used that \mathcal{F} is convex to obtain the set inclusion. Thus, d_0 is feasible for $\phi(x_0)$ and hence

$$\phi(x_0) \leq \langle g, d_0 \rangle = \lambda \langle g, d_1 \rangle + (1 - \lambda) \langle g, d_2 \rangle = \lambda \phi(x_1) + (1 - \lambda) \phi(x_2),$$

which proves that $\phi(x)$ is convex in \mathcal{F}_1 . \square

We are now in position to prove that the criticality measure $\chi(x)$ is Lipschitz continuous on closed and bounded subsets of \mathcal{F} .

THEOREM 3.4 Suppose that AS1, AS2 and AS3b hold. Suppose also that $\nabla_x f(x)$ is Lipschitz continuous on \mathcal{F}_0 with constant κ_{Lg} . Then, there exists a constant $\kappa_{L\chi} > 0$ such that

$$|\chi(x) - \chi(y)| \leq \kappa_{L\chi} \|x - y\| \quad (3.24)$$

for all $x, y \in \mathcal{F}_0$.

Proof. We have from (3.1) that

$$\chi(x) - \chi(y) = \min_{y+d \in \mathcal{F}, \|d\| \leq 1} \langle \nabla_x f(y), d \rangle - \min_{x+d \in \mathcal{F}, \|d\| \leq 1} \langle \nabla_x f(x), d \rangle, \quad (3.25)$$

$$\begin{aligned} &= \min_{y+d \in \mathcal{F}, \|d\| \leq 1} \langle \nabla_x f(y), d \rangle - \min_{y+d \in \mathcal{F}, \|d\| \leq 1} \langle \nabla_x f(x), d \rangle \\ &\quad + \min_{y+d \in \mathcal{F}, \|d\| \leq 1} \langle \nabla_x f(x), d \rangle - \min_{x+d \in \mathcal{F}, \|d\| \leq 1} \langle \nabla_x f(x), d \rangle. \end{aligned} \quad (3.26)$$

Note that the first two terms in (3.26) have the same feasible set but different objectives, while the last two have different feasible sets but the same objective. Consider the difference of the first two terms. Letting

$$\langle \nabla_x f(y), d_y \rangle = \min_{y+d \in \mathcal{F}, \|d\| \leq 1} \langle \nabla_x f(y), d \rangle \quad \text{and} \quad \langle \nabla_x f(x), d_x \rangle = \min_{y+d \in \mathcal{F}, \|d\| \leq 1} \langle \nabla_x f(x), d \rangle,$$

the first difference in (3.26) becomes

$$\begin{aligned} \langle \nabla_x f(y), d_y \rangle - \langle \nabla_x f(x), d_x \rangle &= \langle \nabla_x f(y), d_y - d_x \rangle + \langle \nabla_x f(y) - \nabla_x f(x), d_x \rangle \\ &\leq \langle \nabla_x f(y) - \nabla_x f(x), d_x \rangle \\ &\leq \|\nabla_x f(y) - \nabla_x f(x)\| \cdot \|d_x\| \\ &\leq \kappa_{Lg} \|x - y\|, \end{aligned} \quad (3.27)$$

where to obtain the first inequality above, we used that, by definition of d_y and d_x , d_x is now feasible for the constraints of the problem of which d_y is the solution; the last inequality follows from the assumed Lipschitz continuity of ∇f and from the bound $\|d_x\| \leq 1$.

Consider now the second difference in (3.26) (where we have the same objective but different feasible sets). Employing the last displayed expression on Rockafellar (1970, p. 43), the set $\hat{\mathcal{F}}_0$ in AS2 can be written as

$$\hat{\mathcal{F}}_0 = \mathcal{F}_0 + \delta\mathcal{B},$$

where \mathcal{B} is the open Euclidean unit ball. It is straightforward to show that $\hat{\mathcal{F}}_0 \subseteq \mathcal{F}_1$, where \mathcal{F}_1 is defined by (3.23). Thus, by Lemma 3.3 with $g = \nabla_x f(x)$, ϕ is a proper convex function on $\hat{\mathcal{F}}_0$. This and Rockafellar (1970, Theorem 10.4) now yield that ϕ is Lipschitz continuous (with constant $\kappa_{L\phi}$, say) on any closed and bounded subset of the relative interior of $\hat{\mathcal{F}}_0$, in particular on \mathcal{F}_0 , since $\hat{\mathcal{F}}_0$ is full-dimensional and open and $\mathcal{F}_0 \subseteq \hat{\mathcal{F}}_0$. As a consequence we obtain from (3.26) and (3.27) that

$$\chi(x) - \chi(y) \leq (\kappa_{Lg} + \kappa_{L\phi})\|x - y\|.$$

Since the role of x and y can be interchanged in the above argument, the conclusion of the theorem follows by setting $\kappa_{L\chi} = \kappa_{Lg} + \kappa_{L\phi}$. \square

This theorem provides a generalization of a result already known for the special case where \mathcal{F} is defined by simple bounds and the norm used in the definition of $\chi(x)$ is the infinity norm (see Gratton *et al.*, 2008a, Lemma 4.1).

Next, we prove a first crude upper bound on the length of any model descent step.

LEMMA 3.5 Suppose that AS4 holds and that a given s yields

$$m_k(x_k + s) \leq f(x_k). \quad (3.28)$$

Then

$$\|s\| \leq \frac{3}{\sigma_k} \left(\kappa_B + \sqrt{\sigma_k \|g_k\|} \right). \quad (3.29)$$

Proof. The definition (2.2) and (3.28) give that

$$\langle g_k, s \rangle + \frac{1}{2} \langle s, B_k s \rangle + \frac{1}{3} \sigma_k \|s\|^3 \leq 0.$$

Hence, using the Cauchy–Schwarz inequality and (3.2), we deduce

$$0 \leq \frac{1}{3} \sigma_k \|s\|^3 \leq \|g_k\| \cdot \|s\| + \frac{1}{2} \kappa_B \|s\|^2.$$

This in turn implies that

$$\|s\| \leq \frac{\frac{1}{2} \kappa_B + \sqrt{\frac{1}{4} \kappa_B^2 + \frac{4}{3} \sigma_k \|g_k\|}}{\frac{2}{3} \sigma_k} \leq \frac{\kappa_B + \sqrt{\frac{4}{3} \sigma_k \|g_k\|}}{\frac{2}{3} \sigma_k} \leq \frac{3}{\sigma_k} \left(\kappa_B + \sqrt{\sigma_k \|g_k\|} \right).$$

\square

Using this bound we next verify that Step 1 of Algorithm COCARC is well defined and delivers a suitable generalized Cauchy point.

LEMMA 3.6 Suppose that AS1, AS2 and AS4 hold. Then, for each k with $\chi_k > 0$, the loop between Steps 1.1, 1.2 and 1.3 of Algorithm COCARC is finite and produces a generalized Cauchy point x_k^{GC} satisfying (2.4) and either (2.5) or (2.6).

Proof. Observe first that the generalized Cauchy point resulting from Step 1 must satisfy conditions (2.4), and (2.5) or (2.6), if the loop on j internal to this step terminates finitely. Thus, we only need to show (by contradiction) that this finite termination always occurs. We therefore assume that the loop is infinite and j tends to infinity.

Suppose first that $t_{\max} = \infty$ for all $j \geq 0$. From Lemma 3.5, we know that $\theta(x_k, t_j) = \|x_{k,j} - x_k\|$ is bounded above as a function of j , but yet $t_{j+1} = 2t_j$ and thus t_j tends to infinity. We may then apply (3.16) to deduce that

$$\|P_{\mathcal{T}(x_{k,j})}[-g_k]\| \leq \frac{\theta(x_k, t_j)}{t_j},$$

and thus that

$$\lim_{j \rightarrow \infty} \|P_{\mathcal{T}(x_{k,j})}[-g_k]\| = 0. \quad (3.30)$$

But the same argument that gave (3.14) in Lemma 3.2 implies that, for all $j \geq 0$,

$$-\langle g_k, x_{k,j} - x_k \rangle = |\langle g_k, x_{k,j} - x_k \rangle| = \chi(x_k, \|x_{k,j} - x_k\|).$$

Therefore, Lemma 3.1(iv) provides that $|\langle g_k, x_{k,j} - x_k \rangle|$ is nondecreasing with j and also gives the first inequality below

$$|\langle g_k, x_{k,0} - x_k \rangle| = \chi(x_k, \|x_{k,0} - x_k\|) \geq \min[1, \|x_{k,0} - x_k\|] \chi_k > 0,$$

where the last inequality follows from the fact that x_k is not first-order critical. As a consequence,

$$|\langle g_k, x_{k,j} - x_k \rangle| \geq \min[1, \|x_{k,0} - x_k\|] \chi_k > 0$$

for all $j \geq 0$. Combining this observation with (3.30), we conclude that (2.6) must hold for all j sufficiently large, and the loop inside Step 1 must then be finite, which contradicts our assumption. Thus, our initial supposition on t_{\max} is impossible and t_{\max} must be reset to a finite value. The continuity of the model m_k and of the projection operator $P_{\mathcal{F}}$ then imply, together with (2.7), the existence of an interval I of \mathfrak{R}^+ of nonzero length, possibly nonunique, such that, for all $t \in I$,

$$m_k(P_{\mathcal{F}}[x_k - t g_k]) \leq f(x_k) + \kappa_{\text{ubs}} \langle g_k, P_{\mathcal{F}}[x_k - t g_k] - x_k \rangle$$

and

$$m_k(P_{\mathcal{F}}[x_k - t g_k]) \geq f(x_k) + \kappa_{\text{lbs}} \langle g_k, P_{\mathcal{F}}[x_k - t g_k] - x_k \rangle.$$

But this interval is independent of j and is always contained in $[t_{\min}, t_{\max}]$ by construction, while the length of this latter interval converges to zero when j tends to infinity. Hence, there must exist a finite j such that both (2.4) and (2.5) hold, leading to the desired contradiction. \square

We now derive two finer upper bounds on the length of the generalized Cauchy step, depending on two different criticality measures. These results are inspired by [Cartis et al. \(2011a\)](#), Lemma 2.1).

LEMMA 3.7 Suppose that AS1 and AS2 hold. Then, we have that

$$\|s_k^{\text{GC}}\| \leq \frac{3}{\sigma_k} \max \left[\|B_k\|, (\sigma_k \chi_k)^{\frac{1}{2}}, (\sigma_k^2 \chi_k)^{\frac{1}{3}} \right], \quad (3.31)$$

and

$$\|s_k^{\text{GC}}\| \leq \frac{3}{\sigma_k} \max \left[\|B_k\|, (\sigma_k \pi_k^{\text{GC}})^{\frac{1}{2}} \right], \quad (3.32)$$

Proof. For brevity we omit the index k . From (2.2), (3.14) and the Cauchy–Schwarz inequality,

$$\begin{aligned} m(x^{\text{GC}}) - f(x) &= \langle g, s^{\text{GC}} \rangle + \frac{1}{2} \langle s^{\text{GC}}, B s^{\text{GC}} \rangle + \frac{1}{3} \sigma \|s^{\text{GC}}\|^3 \\ &\geq -\chi(x, \|s^{\text{GC}}\|) - \frac{1}{2} \|s^{\text{GC}}\|^2 \|B\| + \frac{1}{3} \sigma \|s^{\text{GC}}\|^3 \\ &= \left[\frac{1}{9} \sigma \|s^{\text{GC}}\|^3 - \chi(x, \|s^{\text{GC}}\|) \right] + \left[\frac{2}{9} \sigma \|s^{\text{GC}}\|^3 - \frac{1}{2} \|s^{\text{GC}}\|^2 \|B\| \right], \end{aligned} \quad (3.33)$$

Thus, since $m(x^{\text{GC}}) \leq f(x)$, at least one of the bracketed expressions must be negative, i.e. either

$$\|s^{\text{GC}}\| \leq \frac{9}{4} \cdot \frac{\|B\|}{\sigma} \quad (3.34)$$

or

$$\|s^{\text{GC}}\|^3 \leq \frac{9}{\sigma} \chi(x, \|s^{\text{GC}}\|); \quad (3.35)$$

the latter is equivalent to

$$\|s^{\text{GC}}\| \leq 3 \left(\frac{\pi^{\text{GC}}}{\sigma} \right)^{\frac{1}{2}} \quad (3.36)$$

from (3.5) when $\theta = \|s^{\text{GC}}\|$. In the case that $\|s^{\text{GC}}\| \geq 1$, (3.10) then gives that

$$\|s^{\text{GC}}\| \leq 3 \left(\frac{\chi}{\sigma} \right)^{\frac{1}{2}}. \quad (3.37)$$

Conversely, if $\|s^{\text{GC}}\| < 1$, we obtain from (3.11) and (3.35) that

$$\|s^{\text{GC}}\| \leq 3 \left(\frac{\chi}{\sigma} \right)^{\frac{1}{3}}. \quad (3.38)$$

Gathering (3.34), (3.37) and (3.38) we immediately obtain (3.31). Combining (3.34) and (3.36) gives (3.32). \square

Similar results may then be derived for the length of the full step, as we show next.

LEMMA 3.8 Suppose that AS1 and AS2 hold. Then

$$\|s_k\| \leq \frac{3}{\sigma_k} \max \left[\|B_k\|, (\sigma_k \chi_k)^{\frac{1}{2}}, (\sigma_k^2 \chi_k)^{\frac{1}{3}} \right] \quad (3.39)$$

and

$$\|s_k\| \leq \frac{3}{\sigma_k} \max \left[\|B_k\|, \sqrt{\sigma_k \pi_k^{\text{GC}}} \right]. \quad (3.40)$$

Proof. We start by proving (3.39) and

$$\|s_k\| \leq \frac{3}{\sigma_k} \max \left[\|B_k\|, \sqrt{\sigma_k \pi_k^+} \right] \quad (3.41)$$

in a manner identical to that used for (3.31) and (3.32) with s_k replacing s_k^{GC} ; instead of using (3.14) in (3.33) we now employ the inequality $\langle g_k, s_k \rangle \geq -\chi(x_k, \|s_k\|)$, which follows from (3.1). Also, in order to derive the analogues of (3.37) and (3.38), we use (3.12) and (3.13) instead of (3.10) and (3.11), respectively.

If $\|s_k\| \leq \|s_k^{\text{GC}}\|$ then (3.40) immediately follows from (3.32). Otherwise, i.e. if $\|s_k\| > \|s_k^{\text{GC}}\|$ then the nonincreasing nature of $\pi(x_k, \theta)$ gives that $\pi_k^+ \leq \pi_k^{\text{GC}}$. Substituting the latter inequality in (3.41) gives (3.40) in this case. \square

Using the above results we may then derive the equivalent of the well-known Cauchy decrease condition in our constrained case. Again, the exact expression of this condition depends on the criticality measure being considered.

LEMMA 3.9 Suppose that AS1 and AS2 hold. If (2.5) holds and $\|s_k^{\text{GC}}\| \leq 1$ then

$$f(x_k) - m_k(x_k^{\text{GC}}) \geq \kappa_{\text{GC}} \pi_k^{\text{GC}} \min \left[\frac{\pi_k^{\text{GC}}}{1 + \|B_k\|}, \sqrt{\frac{\pi_k^{\text{GC}}}{\sigma_k}} \right], \quad (3.42)$$

where $\kappa_{\text{GC}} \stackrel{\text{def}}{=} \frac{1}{2} \kappa_{\text{ubs}} (1 - \kappa_{\text{lbs}}) \in (0, 1)$. Otherwise, if (2.5) fails and $\|s_k^{\text{GC}}\| \leq 1$, or if $\|s_k^{\text{GC}}\| \geq 1$, then

$$f(x_k) - m_k(x_k^{\text{GC}}) \geq \kappa_{\text{GC}} \chi_k. \quad (3.43)$$

If $\|s_k^{\text{GC}}\| \leq 1$ then

$$f(x_k) - m_k(x_k^{\text{GC}}) \geq \kappa_{\text{GC}} \chi_k \min \left[\frac{\chi_k}{1 + \|B_k\|}, \sqrt{\frac{\pi_k^{\text{GC}}}{\sigma_k}}, 1 \right]. \quad (3.44)$$

In all cases

$$f(x_k) - m_k(x_k^{\text{GC}}) \geq \kappa_{\text{GC}} \chi_k \min \left[\frac{\chi_k}{1 + \|B_k\|}, \sqrt{\frac{\chi_k}{\sigma_k}}, 1 \right]. \quad (3.45)$$

Proof. Again, we omit the index k for brevity. Note that, because of (2.4) and (3.14),

$$f(x) - m(x^{\text{GC}}) \geq \kappa_{\text{ubs}} \left| \langle g, s^{\text{GC}} \rangle \right| = \kappa_{\text{ubs}} \chi(x, \|s^{\text{GC}}\|) = \kappa_{\text{ubs}} \left(x, \|s^{\text{GC}}\| \right) \|s^{\text{GC}}\|. \quad (3.46)$$

Assume first that $\|s^{\text{GC}}\| \geq 1$. Then, using (3.10), we see that

$$f(x) - m(x^{\text{GC}}) \geq \kappa_{\text{ubs}} \chi, \quad (3.47)$$

which gives (3.43) in the case $\|s^{\text{GC}}\| \geq 1$ since $\kappa_{\text{ubs}} > \kappa_{\text{GC}}$. Assume now, for the remainder of the proof, that $\|s^{\text{GC}}\| \leq 1$, which implies, by (3.11), that

$$f(x) - m(x^{\text{GC}}) \geq \kappa_{\text{ubs}} \chi \|s^{\text{GC}}\|, \quad (3.48)$$

and first consider the case where (2.5) holds. Then, from (2.2) and (2.5), the Cauchy–Schwarz inequality, (3.14) and (3.5), we obtain that

$$\|B\| + \frac{2}{3}\sigma \|s^{\text{GC}}\| \geq \frac{2(1 - \kappa_{\text{lbs}})}{\|s^{\text{GC}}\|^2} \left| \langle g, s^{\text{GC}} \rangle \right| = \frac{2(1 - \kappa_{\text{lbs}})}{\|s^{\text{GC}}\|^2} \chi(x, \|s^{\text{GC}}\|) = \frac{2(1 - \kappa_{\text{lbs}})}{\|s^{\text{GC}}\|} \pi^{\text{GC}},$$

and hence that

$$\|s^{\text{GC}}\| \geq \frac{2(1 - \kappa_{\text{lbs}})\pi^{\text{GC}}}{\|B\| + \frac{2}{3}\sigma \|s^{\text{GC}}\|}.$$

Recalling (3.32) we thus deduce that

$$\|s^{\text{GC}}\| \geq \frac{2(1 - \kappa_{\text{lbs}})\pi^{\text{GC}}}{\|B\| + 2 \max \left[\|B\|, \sqrt{\sigma \pi^{\text{GC}}} \right]}.$$

Combining this inequality with (3.46) we obtain that

$$f(x) - m(x^{\text{GC}}) \geq \frac{2}{3} \kappa_{\text{ubs}} (1 - \kappa_{\text{lbs}}) \pi^{\text{GC}} \min \left[\frac{\pi^{\text{GC}}}{1 + \|B\|}, \sqrt{\frac{\pi^{\text{GC}}}{\sigma}} \right],$$

which implies (3.42).

If (2.5) does not hold (and $\|s_k^{\text{GC}}\| \leq 1$) then (2.6) must hold. Thus, (3.15) and (2.7) imply that

$$\chi \leq (1 + 2\kappa_{\text{ep}}) \chi(x, \|s^{\text{GC}}\|) \leq 2\chi(x, \|s^{\text{GC}}\|).$$

Substituting this inequality in (3.46) then gives that

$$f(x) - m(x^{\text{GC}}) \geq \frac{1}{2} \kappa_{\text{ubs}} \chi. \quad (3.49)$$

This in turn implies (3.43) for the case when (2.5) fails and $\|s_k^{\text{GC}}\| \leq 1$. The inequality (3.44) results from (3.42) and (3.11) in the case when (2.5) holds and from (3.49) when (2.5) does not hold. Finally, (3.45) follows from combining (3.42) and (3.43) and using (3.11) in the former. \square

We next show that when the iterate x_k is sufficiently noncritical, then iteration k must be very successful and the regularization parameter does not increase.

LEMMA 3.10 Suppose AS1, AS2 and AS4 hold, that $\chi_k > 0$ and that

$$\min \left[\sigma_k, (\sigma_k \chi_k)^{\frac{1}{2}}, (\sigma_k^2 \chi_k)^{\frac{1}{3}} \right] \geq \frac{9(\kappa_H + \kappa_B)}{2(1 - \eta_2)\kappa_{GC}} \stackrel{\text{def}}{=} \kappa_{\text{suc}} > 1, \quad (3.50)$$

where κ_{GC} is defined just after (3.42). Then, iteration k is very successful and

$$\sigma_{k+1} \leq \sigma_k. \quad (3.51)$$

Proof. First, note that the last inequality in (3.50) follows from the facts that $\kappa_H \geq 1$, $\kappa_B \geq 1$ and $\kappa_{GC} \in (0, 1)$. Again, we omit the index k for brevity. The mean-value theorem gives that

$$f(x^+) - m(x^+) = \frac{1}{2} \langle s, [H(\xi) - B]s \rangle - \frac{1}{3} \sigma \|s\|^3$$

for some $\xi \in [x, x^+]$. Hence, using (3.2),

$$f(x^+) - m(x^+) \leq \frac{1}{2} (\kappa_H + \kappa_B) \|s\|^2. \quad (3.52)$$

We also note that (3.50) and AS4 imply that $(\sigma \chi)^{\frac{1}{2}} \geq \|B\|$ and hence, from (3.39), that

$$\|s\| \leq \frac{3}{\sigma} \max \left[(\sigma \chi)^{\frac{1}{2}}, (\sigma^2 \chi)^{\frac{1}{3}} \right] = 3 \max \left[\left(\frac{\chi}{\sigma} \right)^{\frac{1}{2}}, \left(\frac{\chi}{\sigma} \right)^{\frac{1}{3}} \right].$$

Substituting this last bound in (3.52) then gives that

$$f(x^+) - m(x^+) \leq \frac{9(\kappa_H + \kappa_B)}{2} \max \left[\frac{\chi}{\sigma}, \left(\frac{\chi}{\sigma} \right)^{\frac{2}{3}} \right]. \quad (3.53)$$

Assume now that $\|s^{\text{GC}}\| \leq 1$ and (2.6) holds but not (2.5), or that $\|s^{\text{GC}}\| > 1$. Then (2.9) and (3.43) also imply that

$$f(x) - m(x^+) \geq f(x) - m(x^{\text{GC}}) \geq \kappa_{GC} \chi.$$

Thus, using this bound and (3.53),

$$\begin{aligned} 1 - \rho &= \frac{f(x^+) - m(x^+)}{f(x) - m(x^+)} \\ &\leq \frac{9(\kappa_H + \kappa_B)}{2\kappa_{GC}\chi} \max \left[\frac{\chi}{\sigma}, \left(\frac{\chi}{\sigma} \right)^{\frac{2}{3}} \right] \\ &= \frac{9(\kappa_H + \kappa_B)}{2\kappa_{GC}} \max \left[\frac{1}{\sigma}, \frac{1}{(\sigma^2 \chi)^{\frac{1}{3}}} \right] \\ &\leq 1 - \eta_2, \end{aligned} \quad (3.54)$$

where the last inequality results from (3.50). Assume alternatively that $\|s^{\text{GC}}\| \leq 1$ and (2.5) holds. We then deduce from (3.11), (3.50) and (3.2) that

$$\sqrt{\sigma \pi^{\text{GC}}} \geq \sqrt{\sigma \chi} \geq 1 + \|B\|. \quad (3.55)$$

Then (3.40) yields that

$$\|s\| \leq 3\sqrt{\frac{\pi^{\text{GC}}}{\sigma}},$$

which can be substituted in (3.52) to give

$$f(x^+) - m(x^+) \leq \frac{9}{2}(\kappa_{\text{H}} + \kappa_{\text{B}}) \frac{\pi^{\text{GC}}}{\sigma}. \quad (3.56)$$

On the other hand, (2.9), (3.42) and (3.55) also imply

$$f(x) - m(x^+) \geq f(x) - m(x^{\text{GC}}) \geq \kappa_{\text{GC}} \pi^{\text{GC}} \sqrt{\frac{\pi^{\text{GC}}}{\sigma}}.$$

Thus, using this last bound, (2.8), (3.56), (3.11) and (3.50), we obtain that

$$1 - \rho = \frac{f(x^+) - m(x^+)}{f(x) - m(x^+)} \leq \frac{9(\kappa_{\text{H}} + \kappa_{\text{B}})}{2\kappa_{\text{GC}}\sqrt{\sigma\pi^{\text{GC}}}} \leq \frac{9(\kappa_{\text{H}} + \kappa_{\text{B}})}{2\kappa_{\text{GC}}\sqrt{\sigma\chi}} \leq 1 - \eta_2. \quad (3.57)$$

We then conclude from (3.54) and (3.57) that $\rho \geq \eta_2$ whenever (3.50) holds, which means that the iteration is very successful and (3.51) follows. \square

Our next result shows that the regularization parameter must remain bounded above unless a critical point is approached. Note that this result does not depend on the objective's Hessian being Lipschitz continuous.

LEMMA 3.11 Suppose that AS1, AS2 and AS4 hold, and that there is a constant $\epsilon \in (0, 1]$ and an index $j \leq \infty$ such that

$$\chi_k \geq \epsilon \quad (3.58)$$

for all $k = 0, \dots, j$. Then, for all $k \leq j$,

$$\sigma_k \leq \max \left[\sigma_0, \frac{\gamma_2 \kappa_{\text{suc}}^2}{\epsilon} \right] \stackrel{\text{def}}{=} \kappa_{\sigma}, \quad (3.59)$$

where κ_{suc} is defined in (3.50).

Proof. Let us first show that the following implication holds, for any $k = 0, \dots, j$,

$$\sigma_k \geq \frac{\kappa_{\text{suc}}^2}{\epsilon} \implies \sigma_{k+1} \leq \sigma_k. \quad (3.60)$$

The left-hand side of (3.60) implies $\sigma_k \geq \kappa_{\text{suc}}$ because $\kappa_{\text{suc}} > 1$ and $\epsilon < 1$. Moreover, one verifies easily, using (3.58), that it also gives

$$(\sigma_k \chi_k)^{\frac{1}{2}} \geq (\sigma_k \epsilon)^{\frac{1}{2}} = (\kappa_{\text{suc}}^2)^{\frac{1}{2}} = \kappa_{\text{suc}},$$

and

$$(\sigma_k^2 \chi_k)^{\frac{1}{3}} \geq \left(\frac{\kappa_{\text{suc}}^4}{\epsilon} \right)^{\frac{1}{3}} \geq (\kappa_{\text{suc}}^3)^{\frac{1}{3}} = \kappa_{\text{suc}}.$$

Hence, we deduce that the left-hand side of (3.60) implies that (3.50) holds; and so (3.51) follows by Lemma 3.10, which is the right-hand side of the implication (3.60).

Thus, when $\sigma_0 \leq \gamma_2 \kappa_{\text{suc}}^2 / \epsilon$, (3.60) provides $\sigma_k \leq \gamma_2 \kappa_{\text{suc}}^2 / \epsilon$ for all $k \leq j$, where we have introduced the factor γ_2 for the case when σ_k is less than $\kappa_{\text{suc}}^2 / \epsilon$ and iteration k is not very successful. Thus, (3.59) holds. Letting $k = 0$ in (3.60) gives (3.59) when $\sigma_0 > \gamma_2 \kappa_{\text{suc}}^2 / \epsilon$ since $\gamma_2 > 1$. \square

We are now ready to prove our first-order convergence result. We first state it for the case where there are only finitely many successful iterations.

LEMMA 3.12 Suppose that AS1, AS2 and AS4 hold and that there are only finitely many successful iterations. Then, $x_k = x_*$ for all sufficiently large k and x_* is first-order critical.

Proof. Clearly, (3.61) holds if the algorithm terminates finitely, i.e. there exists k such that $\chi_k = 0$ (see Step 1 of COCARC); hence, let us assume that $\chi_k > 0$ for all $k \geq 0$. After the last successful iterate is computed, indexed by say k_0 , the construction of the COCARC algorithm implies that $x_{k_0+1} = x_{k_0+i} \stackrel{\text{def}}{=} x_*$ for all $i \geq 1$. Since all iterations $k \geq k_0 + 1$ are unsuccessful, σ_k increases by at least a fraction γ_1 so that $\sigma_k \rightarrow \infty$ as $k \rightarrow \infty$. If $\chi_{k_0+1} > 0$ then $\chi_k = \chi_{k_0+1} > 0$ for all $k \geq k_0 + 1$ and so $\chi_k \geq \min(\chi_0, \dots, \chi_{k_0+1}) \stackrel{\text{def}}{=} \epsilon > 0$ for all k . Lemma 3.11 with $j = \infty$ implies that σ_k is bounded above for all k and we have reached a contradiction. \square

We conclude this section by showing the desired convergence when the number of successful iterations is infinite. As for trust-region methods this is accomplished by first showing first-order criticality along a subsequence of the iterates.

THEOREM 3.13 Suppose that AS1–AS3a and AS4 hold. Then, we have that

$$\liminf_{k \rightarrow \infty} \chi_k = 0. \quad (3.61)$$

Hence, at least one limit point of the sequence $\{x_k\}$ (if any) is first-order critical.

Proof. Clearly, (3.61) holds if the algorithm terminates finitely, i.e. there exists k such that $\chi_k = 0$ (see Step 1 of COCARC); hence, let us assume that $\chi_k > 0$ for all $k \geq 0$. Furthermore, the conclusion also holds when there are finitely many successful iterations because of Lemma 3.12. Suppose therefore that there are infinitely many successful iterations. Assume also that (3.58) holds for all k (with $j = \infty$). The mechanism of the algorithm then implies that, if iteration k is successful,

$$f(x_k) - f(x_{k+1}) \geq \eta_1 [f(x_k) - m_k(x_k^+)] \geq \eta_1 \kappa_{\text{GC}} \chi_k \min \left[\frac{\chi_k}{1 + \|B_k\|}, \sqrt{\frac{\chi_k}{\sigma_k}}, 1 \right],$$

where we have used (2.9) and (3.45) to obtain the last inequality. The bounds (3.2), (3.58) and (3.59) then yield that

$$f(x_k) - f(x_{k+1}) \geq \eta_1 \kappa_{\text{GC}} \epsilon \min \left[\frac{\epsilon}{1 + \kappa_B}, \sqrt{\frac{\epsilon}{\kappa_\sigma}}, 1 \right] \stackrel{\text{def}}{=} \kappa_\epsilon > 0. \quad (3.62)$$

Summing over all successful iterations from 0 to k we deduce that

$$f(x_0) - f(x_{k+1}) = \sum_{j=0, j \in \mathcal{S}}^k [f(x_j) - f(x_{j+1})] \geq i_k \kappa_\epsilon,$$

where i_k denotes the number of successful iterations up to iteration k . Since i_k tends to infinity by assumption, we obtain that the sequence $\{f(x_k)\}$ tends to minus infinity, which is impossible because f is bounded below on \mathcal{F} due to AS3a and $x_k \in \mathcal{F}$ for all k . Hence, (3.58) cannot hold for all $k < \infty$; since ϵ in (3.58) was arbitrary in $(0, 1]$, (3.61) follows. \square

We finally prove that the conclusion of the last theorem is not restricted to a subsequence but holds for the complete sequence of iterates.

THEOREM 3.14 Suppose that AS1–AS4 hold. Then, we have that

$$\lim_{k \rightarrow \infty} \chi_k = 0, \quad (3.63)$$

and all limit points of the sequence $\{x_k\}$ are first-order critical.

Proof. Clearly, if the algorithm has finite termination, i.e. $\chi_k = 0$ for some k , the conclusion follows. If \mathcal{S} is finite the conclusion also follows, directly from Lemma 3.12. Suppose therefore that there are infinitely many successful iterations and that there exists a subsequence $\{t_i\} \subseteq \mathcal{S}$ such that

$$\chi_{t_i} \geq 2\epsilon \quad (3.64)$$

for some $\epsilon > 0$. From (3.61) we deduce the existence of another subsequence $\{\ell_i\} \subseteq \mathcal{S}$ such that, for all i , ℓ_i is the index of the first successful iteration after iteration t_i such that

$$\chi_k \geq \epsilon \text{ for } t_i \leq k < \ell_i \text{ and } \chi_{\ell_i} \leq \epsilon. \quad (3.65)$$

We then define

$$\mathcal{K} = \{k \in \mathcal{S} \mid t_i \leq k < \ell_i\}. \quad (3.66)$$

Thus, for each $k \in \mathcal{K} \subseteq \mathcal{S}$, we obtain from (3.45) and (3.65) that

$$f(x_k) - f(x_{k+1}) \geq \eta_1 [f(x_k) - m_k(x_k^+)] \geq \eta_1 \kappa_{\text{GC}} \epsilon \min \left[\frac{\epsilon}{1 + \|B_k\|}, \sqrt{\frac{\chi_k}{\sigma_k}}, 1 \right]. \quad (3.67)$$

Because $\{f(x_k)\}$ is monotonically decreasing and bounded below, it must be convergent and we thus deduce from (3.67) that

$$\lim_{k \rightarrow \infty, k \in \mathcal{K}} \frac{\chi_k}{\sigma_k} = 0, \quad (3.68)$$

which in turn implies, in view of (3.65), that

$$\lim_{k \rightarrow \infty, k \in \mathcal{K}} \sigma_k = +\infty. \quad (3.69)$$

As a consequence of this limit, (3.31), (3.2) and (3.65), we see that, for $k \in \mathcal{K}$,

$$\|s_k^{\text{GC}}\| \leq 3 \max \left[\frac{\kappa_{\text{B}}}{\sigma_k}, \left(\frac{\chi_k}{\sigma_k} \right)^{\frac{1}{2}}, \left(\frac{\chi_k}{\sigma_k} \right)^{\frac{2}{3}} \right],$$

and thus $\|s_k^{\text{GC}}\|$ converges to zero along \mathcal{K} . We therefore obtain that

$$\|s_k^{\text{GC}}\| < 1 \quad \text{for all } k \in \mathcal{K} \text{ sufficiently large,} \quad (3.70)$$

which implies that (3.44) is applicable for these k , yielding, in view of (3.2) and (3.65), that, for $k \in \mathcal{K}$ sufficiently large,

$$f(x_k) - f(x_{k+1}) \geq \eta_1 [f(x_k) - m_k(x_k^+)] \geq \eta_1 \kappa_{\text{GC}} \epsilon \min \left[\frac{\epsilon}{1 + \kappa_{\text{B}}}, \sqrt{\frac{\pi_k^{\text{GC}}}{\sigma_k}}, 1 \right].$$

But the convergence of the sequence $\{f(x_k)\}$ implies that the left-hand side of this inequality converges to zero and hence that the minimum in the last right-hand side must be attained by its middle term for $k \in \mathcal{K}$ sufficiently large. We therefore deduce that, for these k ,

$$f(x_k) - f(x_{k+1}) \geq \eta_1 \kappa_{\text{GC}} \epsilon \sqrt{\frac{\pi_k^{\text{GC}}}{\sigma_k}}. \quad (3.71)$$

Returning to the sequence of iterates we see that

$$\|x_{\ell_i} - x_{t_i}\| \leq \sum_{k=t_i, k \in \mathcal{K}}^{\ell_i-1} \|x_k - x_{k+1}\| = \sum_{k=t_i, k \in \mathcal{K}}^{\ell_i-1} \|s_k\| \quad \text{for each } \ell_i \text{ and } t_i. \quad (3.72)$$

Recall now the upper bound (3.40) on $\|s_k\|$, $k \geq 0$. It follows from (3.11) that $\pi_k^{\text{GC}} \geq \chi_k \geq \epsilon$, so that (3.69) implies $\sqrt{\sigma_k \pi_k^{\text{GC}}} \geq \kappa_{\text{B}}$ for all $k \in \mathcal{K}$ sufficiently large. Hence, (3.2) and (3.40) ensure the first inequality below,

$$\|s_k\| \leq 3 \sqrt{\frac{\pi_k^{\text{GC}}}{\sigma_k}} \leq \frac{3}{\eta_1 \kappa_{\text{GC}} \epsilon} [f(x_k) - f(x_{k+1})] \quad \text{for } k \in \mathcal{K} \text{ sufficiently large,}$$

where the second inequality follows from (3.71). This last bound can then be used in (3.72) to obtain

$$\|x_{\ell_i} - x_{t_i}\| \leq \frac{3}{\eta_1 \kappa_{\text{GC}} \epsilon} \sum_{k=t_i, k \in \mathcal{K}}^{\ell_i-1} [f(x_k) - f(x_{k+1})] \leq \frac{3}{\eta_1 \kappa_{\text{GC}} \epsilon} [f(x_{t_i}) - f(x_{\ell_i})]$$

for all t_i and ℓ_i sufficiently large. Since $\{f(x_k)\}$ is convergent, the right-hand side of this inequality tends to zero as i tends to infinity. Hence, $\|x_{\ell_i} - x_{t_i}\|$ converges to zero with i , and, by Theorem 3.4, so does $|\chi_{\ell_i} - \chi_{t_i}|$. But this is impossible since (3.64) and (3.65) imply $|\chi_{\ell_i} - \chi_{t_i}| \geq \chi_{t_i} - \chi_{\ell_i} \geq \epsilon$. Hence, no subsequence can exist such that (3.64) holds and the proof is complete. \square

Assumption AS3b in the above theorem is only mildly restrictive and is satisfied if for instance, the feasible set \mathcal{F} itself is bounded, or if the constrained level set of the objective function, $\{x \in \mathcal{F} | f(x) \leq f(x_0)\}$, is bounded. Note also that AS3b would not be required in Theorem 3.14 provided $\chi(x)$ is uniformly continuous on the sequence of iterates.

4. Worst-case function-evaluation complexity

This section is devoted to worst-case function-evaluation complexity bounds; that is bounds on the number of objective-function or gradient evaluations needed to achieve first-order convergence to prescribed accuracy. Despite the obvious observation that such an analysis does not cover the total computational cost of solving a problem, this type of complexity result is of special interest for nonlinear optimization because there are many examples where the cost of these evaluations completely dwarfs that of the other computations inside the algorithm itself.

Note that the construction of the COCARC basic framework implies that the total number of COCARC iterations is the same as the number of objective-function evaluations as we also need to evaluate f on unsuccessful iterations in order to be able to compute ρ_k in (2.8); the number of successful COCARC iterations is the same as the gradient-evaluation count.

Firstly, let us give a generic worst-case result regarding the number of unsuccessful COCARC iterations, namely iterations i with $\rho_i < \eta_1$, that occur up to any given iteration. Given any $j \geq 0$, denote the iteration index sets

$$\mathcal{S}_j \stackrel{\text{def}}{=} \{k \leq j : k \in \mathcal{S}\} \quad \text{and} \quad \mathcal{U}_j \stackrel{\text{def}}{=} \{i \leq j : i \text{ unsuccessful}\}, \quad (4.1)$$

which form a partition of $\{0, \dots, j\}$. Let $|\mathcal{S}_j|$ and $|\mathcal{U}_j|$ denote their respective cardinalities. Concerning σ_k we may require that on each very successful iteration $k \in \mathcal{S}$, i.e. $\rho_k \geq \eta_2$, σ_{k+1} is chosen such that

$$\sigma_{k+1} \geq \gamma_3 \sigma_k \quad \text{for some } \gamma_3 \in (0, 1]. \quad (4.2)$$

Note that (4.2) allows $\{\sigma_k\}$ to converge to zero on very successful iterations (but no faster than $\{\gamma_3^k\}$). A stronger condition on σ_k is

$$\sigma_k \geq \sigma_{\min}, \quad k \geq 0, \quad (4.3)$$

for some $\sigma_{\min} > 0$. The conditions (4.2) and (4.3) will be employed in the complexity bounds for COCARC and a second-order variant, respectively.

THEOREM 4.1 For any fixed $j \geq 0$, let \mathcal{S}_j and \mathcal{U}_j be defined in (4.1). Assume that (4.2) holds and let $\bar{\sigma} > 0$ be such that

$$\sigma_k \leq \bar{\sigma} \quad \text{for all } k \leq j. \quad (4.4)$$

Then,

$$|\mathcal{U}_j| \leq \left\lceil -\frac{\log \gamma_3}{\log \gamma_1} |\mathcal{S}_j| + \frac{1}{\log \gamma_1} \log \left(\frac{\bar{\sigma}}{\sigma_0} \right) \right\rceil. \quad (4.5)$$

In particular, if σ_k satisfies (4.3), then it also achieves (4.2) with $\gamma_3 = \sigma_{\min} / \bar{\sigma}$, and we have that

$$|\mathcal{U}_j| \leq \left\lceil (|\mathcal{S}_j| + 1) \frac{1}{\log \gamma_1} \log \left(\frac{\bar{\sigma}}{\sigma_{\min}} \right) \right\rceil. \quad (4.6)$$

Proof. The proof follows identically to that of *Cartis et al. (2011b, Theorem 2.1)*. \square

4.1 Function-evaluation complexity for COCARC algorithm

We first consider the function- (and gradient-) evaluation complexity of a variant—COCARC $_{\epsilon}$ —of the COCARC algorithm itself, only differing by the introduction of an approximate termination rule. More specifically, we replace the criticality check in Step 1 of COCARC by the test $\chi_k \leq \epsilon$ (where ϵ is a user-supplied threshold) and terminate if this inequality holds. The results presented for this algorithm are inspired by complexity results for trust-region algorithms (see Gratton *et al.*, 2008a,b) and for the adaptive cubic regularization algorithm (see Cartis *et al.*, 2011b).

THEOREM 4.2 Suppose that AS1–AS3a, AS4 and (4.2) hold and that the approximate criticality threshold ϵ is small enough to ensure

$$\epsilon \leq \min \left[1, \frac{\gamma_2 \kappa_{\text{suc}}^2}{\sigma_0} \right], \quad (4.7)$$

where κ_{suc} is defined in (3.50). Assuming $\gamma_0 > \epsilon$ there exists a constant $\kappa_{\text{df}} \in (0, 1)$ such that

$$f(x_k) - f(x_{k+1}) \geq \kappa_{\text{df}} \epsilon^2 \quad (4.8)$$

for all $k \in \mathcal{S}$ before Algorithm COCARC $_{\epsilon}$ terminates, namely, until it generates a first iterate, say x_{j_1} , such that $\chi_{j_1+1} \leq \epsilon$. As a consequence this algorithm needs at most

$$\lceil \kappa_{\mathcal{S}} \epsilon^{-2} \rceil \quad (4.9)$$

successful iterations and evaluations of the objective's gradient $\nabla_x f$ to ensure $\chi_{j_1+1} \leq \epsilon$, and furthermore,

$$j_1 \leq \lceil \kappa_* \epsilon^{-2} \rceil \stackrel{\text{def}}{=} J_1,$$

so that the algorithm takes at most J_1 iterations and objective-function evaluations to terminate with $\chi_{j_1+1} \leq \epsilon$, where

$$\kappa_{\mathcal{S}} \stackrel{\text{def}}{=} \frac{f(x_0) - f_{\text{low}}}{\kappa_{\text{df}}} \quad \text{and} \quad \kappa_* \stackrel{\text{def}}{=} \left(1 - \frac{\log \gamma_3}{\log \gamma_1} \right) \kappa_{\mathcal{S}} + \frac{\gamma_2 \kappa_{\text{suc}}^2}{\sigma_0 \log \gamma_1}.$$

Proof. From the definition of the $(j_1 + 1)$ th iteration we must have $\chi_k > \epsilon$ for all $k \leq j_1$. This, (4.7) and (3.59) imply that

$$\sigma_k \leq \frac{\gamma_2 \kappa_{\text{suc}}^2}{\epsilon} \quad \text{for all } k \leq j_1. \quad (4.10)$$

We may now use the same reasoning as in the proof of Theorem 3.13 and employ (3.62) and (4.10) to deduce that

$$\begin{aligned} f(x_k) - f(x_{k+1}) &\geq \eta_1 \kappa_{\text{GC}} \epsilon \min \left[\frac{\epsilon}{1 + \kappa_{\text{B}}}, \sqrt{\frac{\epsilon}{\gamma_2 \kappa_{\text{suc}}^2 / \epsilon}}, 1 \right] \\ &\geq \eta_1 \kappa_{\text{GC}} \min \left[\frac{1}{1 + \kappa_{\text{H}}}, \frac{1}{\kappa_{\text{suc}} \sqrt{\gamma_2}} \right] \epsilon^2 \quad \text{for all } k \in \mathcal{S}_{j_1}, \end{aligned}$$

where we have used (4.7), namely $\epsilon \leq 1$, to derive the last inequality. This gives (4.8) with

$$\kappa_{\text{df}} \stackrel{\text{def}}{=} \eta_1 \kappa_{\text{GC}} \min \left[\frac{1}{1 + \kappa_{\text{H}}}, \frac{1}{\kappa_{\text{suc}} \sqrt{\gamma_2}} \right].$$

The bound (4.8) and the fact that f does not change on unsuccessful iterations imply

$$f(x_0) - f(x_{j_1+1}) = \sum_{k=0, k \in \mathcal{S}}^{j_1} (f(x_k) - f(x_{k+1})) \geq |\mathcal{S}_{j_1}| \kappa_{\text{df}} \epsilon^2,$$

which, due to AS3a, further gives

$$|\mathcal{S}_{j_1}| \leq \frac{f(x_0) - f_{\text{low}}}{\kappa_{\text{df}}} \epsilon^{-2}. \quad (4.11)$$

This immediately provides (4.9) since $|\mathcal{S}_{j_1}|$ must be an integer. Finally, to bound the total number of iterations up to j_1 , recall (4.2) and employ the upper bound on σ_k given in (4.10) as $\bar{\sigma}$ in (4.5) to deduce

$$|\mathcal{U}_{j_1}| \leq \left\lceil -\frac{\log \gamma_3}{\log \gamma_1} |\mathcal{S}_{j_1}| + \frac{1}{\log \gamma_1} \log \left(\frac{\gamma_2 \kappa_{\text{suc}}^2}{\epsilon \sigma_0} \right) \right\rceil.$$

This, the bound (4.9) on $|\mathcal{S}_{j_1}|$ and the inequality $\log(\gamma_2 \kappa_{\text{suc}}^2 / (\epsilon \sigma_0)) \leq (\gamma_2 \kappa_{\text{suc}}^2 / (\epsilon \sigma_0))$ now imply

$$j_1 = |\mathcal{S}_{j_1}| + |\mathcal{U}_{j_1}| \leq \left\lceil \left(1 - \frac{\log \gamma_3}{\log \gamma_1} \right) \kappa_{\mathcal{S}} \epsilon^{-2} + \frac{\gamma_2 \kappa_{\text{suc}}^2}{\sigma_0 \log \gamma_1} \epsilon^{-1} \right\rceil.$$

The bound on j_1 now follows by using $\epsilon \leq 1$. \square

Because Algorithm COCARC_ϵ does not exploit more than first-order information (via the Cauchy point definition), the above upper bound is, as expected, of the same order in ϵ as that obtained by Nesterov (2004, p. 29), and by Vavasis (1993), for the steepest descent method.

4.2 An $\mathcal{O}(\epsilon^{-3/2})$ function-evaluation complexity bound

We now discuss a variant— COCARC-S —of the COCARC algorithm for which an interesting worst-case function- (and derivatives-) evaluation complexity result can be shown. Algorithm COCARC-S uses the user-supplied first-order accuracy threshold, $\epsilon > 0$. It differs from the basic COCARC framework in that stronger conditions are imposed on the step.

Let us first mention some assumptions on the true and approximate Hessian of the objective that will be required at various points in this section.

AS5. The Hessian $H(x_k)$ is well approximated by B_k , in the sense that there exists a constant $\kappa_{\text{BH}} > 0$ such that, for all k ,

$$\|[B_k - H(x_k)]s_k\| \leq \kappa_{\text{BH}} \|s_k\|^2.$$

AS6. The Hessian of the objective function is ‘weakly’ uniformly Lipschitz continuous on the segments $[x_k, x_k + s_k]$, in the sense that there exists a constant $\kappa_{\text{LH}} \geq 0$ such that, for all k and all $y \in [x_k, x_k + s_k]$,

$$\|[H(y) - H(x_k)]s_k\| \leq \kappa_{\text{LH}} \|s_k\|^2.$$

AS5 and AS6 are acceptable assumptions essentially corresponding to the cases analysed in Nesterov & Polyak (2006) and Cartis *et al.* (2011b) for the unconstrained case, the only differences being that the first authors assume $B_k = H(x_k)$ instead of the weaker AS5.

4.2.1 *A termination condition for the model subproblem.* The conditions on the step in COCARC-S may require the (approximate) constrained model minimization to be performed to higher accuracy than that provided by the Cauchy point. A common way to achieve this is to impose an appropriate termination condition for the inner iterations that perform the constrained model minimization as follows.

AS7: For all k the step s_k solves the subproblem

$$\min_{s \in \mathbb{R}^n, x_k + s \in \mathcal{F}} m_k(x_k + s) \quad (4.12)$$

accurately enough to ensure that

$$\chi_k^m(x_k^+) \leq \min(\kappa_{\text{stop}}, \|s_k\|)\chi_k, \quad (4.13)$$

where $\kappa_{\text{stop}} \in [0, 1)$ is a constant and where

$$\chi_k^m(x) \stackrel{\text{def}}{=} \left| \min_{x+d \in \mathcal{F}, \|d\| \leq 1} \langle \nabla_s m_k(x), d \rangle \right|. \quad (4.14)$$

Note that $\chi_k^m(x_k) = \chi_k$. The inequality (4.13) is an adequate stopping condition for the subproblem solution since $\chi_k^m(x_k^*)$ is equal to zero if x_k^* is a local minimizer of (4.12). It is the constrained analogue of the ‘s-stopping rule’ of [Cartis et al. \(2011b\)](#). Note that though ensuring AS7 may be NP-hard computationally, it does not require any additional objective-function or gradient evaluations, and as such, it will not worsen the global complexity bound for COCARC-S, which counts these evaluations.

An important consequence of AS5–AS7 is that they allow us to deduce the following crucial relation between the local optimality measure and the step.

LEMMA 4.3 i) Suppose that AS1–AS2 and AS5–AS6 hold. Then

$$\sigma_k \leq \max \left[\sigma_0, \frac{3}{2} \gamma_2 (\kappa_{\text{BH}} + \kappa_{\text{LH}}) \right] \stackrel{\text{def}}{=} \sigma_{\text{max}} \quad \text{for all } k \geq 0. \quad (4.15)$$

ii) Suppose that AS1–AS7 hold. Then

$$\|s_k\| \geq \kappa_s \sqrt{\chi_{k+1}} \quad \text{for all } k \in \mathcal{S}, \quad (4.16)$$

for some constant $\kappa_s \in (0, 1)$ independent of k , where χ_k is defined just after (3.1).

Proof. (i) The proof of (4.15) follows identically to that of [Cartis et al. \(2011a, Lemma 5.2\)](#), as the mechanism for updating σ_k and for deciding the success or otherwise of iteration k are identical in the COCARC and the (unconstrained) ARC frameworks.

(ii) Since $k \in \mathcal{S}$ and by definition of the trial point, we have $x_{k+1} = x_k^+ = x_k + s_k$, and hence by (3.1), $\chi_{k+1} = \chi(x_k^+)$. Again, let us drop the index k for the proof, define $\chi^+ \stackrel{\text{def}}{=} \chi(x_k^+)$ and $g^+ \stackrel{\text{def}}{=} g(x_k^+)$, and derive by Taylor expansion of g^+ ,

$$\begin{aligned}
\|g^+ - \nabla_s m(x_k^+)\| &= \left\| g + \int_0^1 H(x + ts)s \, dt - g - [B - H(x)]s - H(x)s - \sigma \|s\|s \right\| \\
&\leq \left\| \int_0^1 [H(x + ts) - H(x)]s \, dt \right\| + (\kappa_{\text{BH}} + \sigma)\|s\|^2 \\
&\leq \int_0^1 \| [H(x + ts) - H(x)]s \| \, dt + (\kappa_{\text{BH}} + \sigma)\|s\|^2 \\
&\leq (\kappa_{\text{LH}} + \kappa_{\text{BH}} + \sigma)\|s\|^2, \\
&\leq (\kappa_{\text{LH}} + \kappa_{\text{BH}} + \sigma_{\text{max}})\|s\|^2,
\end{aligned} \tag{4.17}$$

where we have used (2.2), AS5, AS6, the triangular inequality and (4.15). Assume first that

$$\|s\| \geq \sqrt{\frac{\chi^+}{2(\kappa_{\text{LH}} + \kappa_{\text{BH}} + \sigma_{\text{max}})}}. \tag{4.18}$$

In this case (4.16) follows with $\kappa_s = \sqrt{\frac{1}{2(\kappa_{\text{LH}} + \kappa_{\text{BH}} + \sigma_{\text{max}})}}$, as desired. Assume therefore that (4.18) fails and observe that

$$\chi^+ \stackrel{\text{def}}{=} |\langle g^+, d^+ \rangle| = -\langle g^+, d^+ \rangle \leq |\langle g^+ - \nabla_s m(x^+), d^+ \rangle| + |\langle \nabla_s m(x^+), d^+ \rangle|, \tag{4.19}$$

where the first equality defines the vector d^+ with

$$\|d^+\| \leq 1. \tag{4.20}$$

But, using the Cauchy–Schwarz inequality, (4.20), (4.17), the failure of (4.18) and the first part of (4.19) successively, we obtain

$$\begin{aligned}
\langle \nabla_s m(x^+), d^+ \rangle - \langle g^+, d^+ \rangle &\leq |\langle g^+, d^+ \rangle - \langle \nabla_s m(x^+), d^+ \rangle| \\
&\leq \|g^+ - \nabla_s m(x^+)\| \\
&\leq (\kappa_{\text{LH}} + \kappa_{\text{BH}} + \sigma_{\text{max}})\|s\|^2 \\
&\leq \frac{1}{2}\chi^+ \\
&= -\frac{1}{2}\langle g^+, d^+ \rangle,
\end{aligned}$$

which in turn ensures that

$$\langle \nabla_s m(x^+), d^+ \rangle \leq \frac{1}{2}\langle g^+, d^+ \rangle < 0.$$

Moreover, $x^+ + d^+ \in \mathcal{F}$ by definition of χ^+ , and hence, using (4.20) and (4.14),

$$|\langle \nabla_s m(x^+), d^+ \rangle| \leq \chi^m(x^+). \tag{4.21}$$

We may then substitute this bound in (4.19) and use the Cauchy–Schwarz inequality and (4.20) again to deduce that

$$\chi^+ \leq \|g^+ - \nabla_s m(x^+)\| + \chi^m(x^+) \leq \|g^+ - \nabla_s m(x^+)\| + \min(\kappa_{\text{stop}}, \|s\|)\chi, \tag{4.22}$$

where the last inequality results from (4.13).

We now observe that both x and x^+ belong to \mathcal{F}_0 , where \mathcal{F}_0 is defined in AS1. Moreover, the first inequality in (3.2) provides that $\nabla_x f(x)$ is Lipschitz continuous on \mathcal{F}_0 , with constant $\kappa_{Lg} = \kappa_H$. Thus, Theorem 3.4 applies, ensuring that $\chi(x)$ is Lipschitz continuous on \mathcal{F}_0 , with Lipschitz constant $\kappa_{L\chi}$; it follows from (3.24) applied to x and x^+ that

$$\chi \leq \kappa_{L\chi} \|x - x^+\| + \chi^+ = \kappa_{L\chi} \|s\| + \chi^+, \quad (4.23)$$

which substituted in (4.22), gives

$$\chi^+ \leq \|g^+ - \nabla_s m(x^+)\| + \min(\kappa_{\text{stop}}, \|s\|) [\kappa_{L\chi} \|s\| + \chi^+] \leq \|g^+ - \nabla_s m(x^+)\| + \kappa_{L\chi} \|s\|^2 + \kappa_{\text{stop}} \chi^+,$$

where the second inequality follows by employing $\min(\kappa_{\text{stop}}, \|s\|) \leq \|s\|$ and $\min(\kappa_{\text{stop}}, \|s\|) \leq \kappa_{\text{stop}}$, respectively. Now substituting (4.17) into the last displayed inequality, we obtain

$$\chi^+ \leq (\kappa_{LH} + \kappa_{BH} + \sigma_{\max}) \|s\|^2 + \kappa_{L\chi} \|s\|^2 + \kappa_{\text{stop}} \chi^+,$$

which further gives

$$(1 - \kappa_{\text{stop}}) \chi^+ \leq (\kappa_{LH} + \kappa_{L\chi} + \kappa_{BH} + \sigma_{\max}) \|s\|^2.$$

Therefore, since $\kappa_{\text{stop}} \in (0, 1)$, we deduce

$$\|s\| \geq \sqrt{\frac{(1 - \kappa_{\text{stop}}) \chi^+}{\kappa_{LH} + \kappa_{L\chi} + \kappa_{BH} + \sigma_{\max}}},$$

which gives (4.16) with

$$\kappa_s = \sqrt{\frac{1 - \kappa_{\text{stop}}}{\kappa_{LH} + \kappa_{L\chi} + \kappa_{BH} + \sigma_{\max}}}. \quad (4.24)$$

□

4.2.2 Ensuring the model decrease. Similarly to the unconstrained case presented in [Cartis et al. \(2011b\)](#), AS7 is unfortunately not sufficient to obtain the desired complexity result; in particular, this may not ensure a model decrease of the form

$$m_k(x_k) - m_k(x_k^+) \geq \kappa_{\text{red}} \sigma_k \|s_k\|^3, \quad (4.25)$$

for some constant $\kappa_{\text{red}} > 0$, independent of k , where $m_k(x_k) = f(x_k)$. For x_k^+ to be an acceptable trial point, one also needs to verify that a cheap but too small model improvement cannot be obtained from x_k^+ . In the unconstrained case this was expressed by the requirement that the trial point is a stationary point of the model at least in some subspace and that the step provides a descent direction. [To see why these conditions imply a decrease of type (4.25) in the unconstrained case, see [Cartis et al., 2011a](#), Lemma 3.3.] An even milder form of the former condition can be easily imposed in the constrained case too, by requiring that the step s_k satisfies

$$\langle \nabla_s m_k(x_k^+), s_k \rangle \leq 0, \quad (4.26)$$

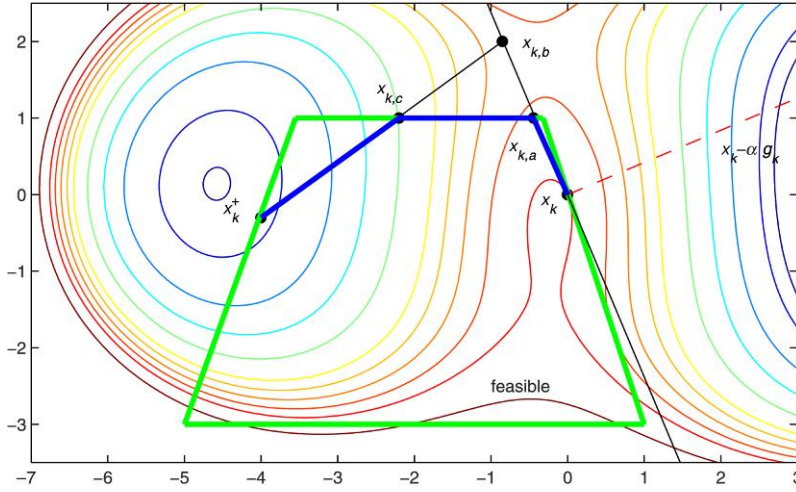


FIG. 1. An illustration when (4.27) fails for the cubic model $m(x, y) = -x - \frac{42}{100}y - \frac{3}{10}x^2 - \frac{1}{10}y^3 + \frac{1}{3}[x^2 + y^2]^{3/2}$ at the iterate $x_k = (0, 0)^T$; the feasible set \mathcal{F} is the polytope with vertices $(1, -5)^T$, $(-\frac{32}{100}, 1)^T$, $(-\frac{355}{100}, 1)^T$ and $(-\frac{510}{100}, -5)^T$. The path \mathcal{P}_k defined in (4.47) that satisfies AS8 is also represented.

which expresses the reasonable requirement that the step size along s_k does not exceed that corresponding to the minimum of the model $m_k(x_k + \tau s_k)$ for $\tau > 0$. It is, for instance, satisfied if

$$1 \in \underset{\tau \geq 0, x_k + \tau s_k \in \mathcal{F}}{\operatorname{argmin}} m_k(x_k + \tau s_k).$$

Note that (4.26) also holds at a local minimizer. Lemma 4.4 below shows that (4.25) is indeed satisfied when (4.26) holds, provided the step s_k is descent or the model is convex.

However, at variance with the unconstrained case, there is no longer any guarantee that the step s_k provides a descent direction in the presence of negative curvature, i.e. that $\langle \nabla_s m_k(x_k), s_k \rangle \leq 0$ when $\langle s_k, B_k s_k \rangle < 0$; recall that $\nabla_s m_k(x_k) = g_k$. Figure 1 illustrates the latter situation; namely, the contours of a particular model $m_k(x_k + s)$ are plotted, as well as a polyhedral feasible set \mathcal{F} , the steepest descent direction from x_k and the hyperplane orthogonal to it, i.e. $\langle \nabla_s m_k(x_k), s \rangle = 0$. Note that all acceptable feasible directions from x_k (pointing towards the feasible local model minimizer) are ascent locally, as the (only) feasible local model minimizer lies on the ‘wrong side’ of the ‘mountain’, in a direction such that (4.27) fails. However, in this unsatisfactory situation, there may be a piecewise-linear feasible descent path (towards the local model minimizer) that goes around the mountain, taking us downhill at each step; see AS8 and the path determined by $\{x_k, x_{k,a}, x_{k,c}, x_k^+\}$ in Fig. 1. In the latter case we will show that the bound (4.25) holds, provided the local path to the trial point x_k^+ contains a uniformly bounded number of descent line segments. Let us now make these illustrations mathematically precise. We begin by considering the easy case.

LEMMA 4.4 Suppose that (4.26) holds and that

$$\langle \nabla_s m_k(x_k), s_k \rangle \leq 0 \quad \text{or} \quad \langle s_k, B_k s_k \rangle \geq 0. \tag{4.27}$$

Then,

$$m_k(x_k) - m_k(x_k^+) \geq \frac{1}{6} \sigma_k \|s_k\|^3. \quad (4.28)$$

Proof. (Dropping the index k again.) Condition (4.26) is equivalent to

$$\langle g, s \rangle + \langle s, Bs \rangle + \sigma \|s\|^3 \leq 0. \quad (4.29)$$

If $\langle s, Bs \rangle \geq 0$, we substitute $\langle g, s \rangle$ from this inequality into (2.2) and deduce that

$$m(x^+) - f(x) = \langle g, s \rangle + \frac{1}{2} \langle s, Bs \rangle + \frac{1}{3} \sigma \|s\|^3 \leq -\frac{1}{2} \langle s, Bs \rangle - \frac{2}{3} \sigma \|s\|^3,$$

which then implies (4.28). If, on the other hand, $\langle s, Bs \rangle < 0$, then we substitute the inequality on $\langle s, Bs \rangle$ resulting from (4.29) into (2.2) and obtain that

$$m(x^+) - f(x) = \langle g, s \rangle + \frac{1}{2} \langle s, Bs \rangle + \frac{1}{3} \sigma \|s\|^3 \leq \frac{1}{2} \langle g, s \rangle - \frac{1}{6} \sigma \|s\|^3$$

from which (4.28) again follows because of (4.27). \square

Note that the following implication follows from (4.29) and $g_k = \nabla_s m_k(x_k)$,

$$(4.26) \quad \text{and} \quad \langle s_k, B_k s_k \rangle \geq 0 \quad \implies \quad \langle \nabla_s m_k(x_k), s_k \rangle \leq 0. \quad (4.30)$$

As we already mentioned, ensuring (4.25) is more complicated when (4.27) fails, namely, when the step is ascent (at x_k) rather than descent and of negative curvature. Our requirement on the trial point is then essentially that it can be computed by a uniformly bounded sequence of (possibly incomplete) line minimizations starting from x_k . More formally, we assume that there exists an integer $\bar{\ell} > 0$ and, for each k such that (4.27) fails, there exist feasible points $\{x_{k,i}\}_{i=0}^{\ell_k}$ with $0 < \ell_k \leq \bar{\ell}$, $x_{k,0} = x_k$ and $x_{k,\ell_k} = x_k^+$, such that, for $i = 1, \dots, \ell_k$,

$$m_k(x_{k,i}) \leq m_k(x_{k,i-1}), \quad \langle \nabla_s m_k(x_{k,i-1}), x_{k,i} - x_{k,i-1} \rangle \leq 0 \quad \text{and} \quad \langle \nabla_s m_k(x_{k,i}), x_{k,i} - x_{k,i-1} \rangle \leq 0. \quad (4.31)$$

Note that these inequalities hold in particular if x_k^+ is the first minimizer of the model along the piecewise linear path

$$\mathcal{P}_k \stackrel{\text{def}}{=} \bigcup_{i=1}^{\ell_k} [x_{k,i-1}, x_{k,i}];$$

such a trial point exists since f is continuous and the path \mathcal{P}_k is compact. The conditions (4.31) subsume the case addressed in Lemma 4.4 when (4.27) holds because one may then choose $\ell_k = 1$ and (4.31) then implies both (4.26) and (4.27); recall also (4.30). We can therefore comprehensively summarize all these requirements in the following assumption.

AS8. For all k the step s_k is such that (4.31) holds for some $\{x_{k,i}\}_{i=0}^{\ell_k} \subset \mathcal{F}$ with $0 < \ell_k \leq \bar{\ell}$, $x_{k,0} = x_k$ and $x_{k,\ell_k} = x_k^+$.

Observe that we have not used global constrained optimization anywhere in the requirements imposed on the step s_k .

Using AS8 we may now obtain the essential lower bound on the model reduction. First, we give a useful technical lemma.

LEMMA 4.5 Suppose that there exist steps $s_{k,\circ}$ and $s_{k,\bullet}$ and points $x_{k,\circ} = x_k + s_{k,\circ}$ and $x_{k,\bullet} = x_k + s_{k,\bullet}$ such that, for some $\kappa \in (0, 1]$,

$$m_k(x_{k,\circ}) \leq m_k(x_k) - \kappa\sigma_k \|s_{k,\circ}\|^3, \quad (4.32)$$

$$m_k(x_{k,\bullet}) \leq m_k(x_{k,\circ}), \quad (4.33)$$

$$\langle \nabla_s m_k(x_{k,\bullet}), x_{k,\bullet} - x_{k,\circ} \rangle \leq 0, \quad (4.34)$$

and

$$\langle \nabla_s m_k(x_{k,\circ}), x_{k,\bullet} - x_{k,\circ} \rangle \leq 0. \quad (4.35)$$

Then,

$$m_k(x_k) - m_k(x_{k,\bullet}) \geq \kappa_{\text{lm}} \kappa \sigma_k \|s_{k,\bullet}\|^3 \quad (4.36)$$

for some constant $\kappa_{\text{lm}} \in (0, 1)$ independent of k and κ .

Proof. (Dropping the index k again.) Suppose first that, for some $\alpha \in (0, 1)$,

$$\|s_\circ\| \geq \alpha \|s_\bullet\|. \quad (4.37)$$

Then, (4.32) and (4.33) give that

$$m(x) - m(x_\bullet) = m(x) - m(x_\circ) + m(x_\circ) - m(x_\bullet) \geq \kappa\sigma \|s_\circ\|^3 \geq \kappa\sigma\alpha^3 \|s_\bullet\|^3. \quad (4.38)$$

Assume now that (4.37) fails; that is

$$\|s_\circ\| < \alpha \|s_\bullet\|. \quad (4.39)$$

We have that

$$f(x) + \langle g, s_\circ \rangle + \frac{1}{2} \langle s_\circ, Bs_\circ \rangle = m(x_\circ) - \frac{1}{3} \sigma \|s_\circ\|^3. \quad (4.40)$$

Using this identity we now see that

$$\begin{aligned} m(x_\bullet) &= f(x) + \langle g, s_\circ \rangle + \frac{1}{2} \langle s_\circ, Bs_\circ \rangle + \langle g + Bs_\circ, s_\bullet - s_\circ \rangle + \frac{1}{2} \langle s_\bullet - s_\circ, B(s_\bullet - s_\circ) \rangle + \frac{1}{3} \sigma \|s_\bullet\|^3 \\ &= m(x_\circ) + \langle g + Bs_\circ, s_\bullet - s_\circ \rangle + \frac{1}{2} \langle s_\bullet - s_\circ, B(s_\bullet - s_\circ) \rangle + \frac{1}{3} \sigma \|s_\bullet\|^3 - \frac{1}{3} \sigma \|s_\circ\|^3. \end{aligned} \quad (4.41)$$

Moreover, (4.34) yields that

$$\begin{aligned} 0 &\geq \langle g + Bs_\circ, s_\bullet - s_\circ \rangle + \sigma \|s_\bullet\| \langle s_\bullet, s_\bullet - s_\circ \rangle \\ &= \langle g + Bs_\circ, s_\bullet - s_\circ \rangle + \langle s_\bullet - s_\circ, B(s_\bullet - s_\circ) \rangle + \sigma \|s_\bullet\| \langle s_\bullet, s_\bullet - s_\circ \rangle, \end{aligned}$$

and thus, (4.41) becomes

$$m(x_{\bullet}) \leq m(x_{\circ}) + \frac{1}{2} \langle g + Bs_{\circ}, s_{\bullet} - s_{\circ} \rangle - \frac{1}{2} \sigma \|s_{\bullet}\| \langle s_{\bullet}, s_{\bullet} - s_{\circ} \rangle + \frac{1}{3} \sigma \|s_{\bullet}\|^3 - \frac{1}{3} \sigma \|s_{\circ}\|^3. \quad (4.42)$$

But we may also use (4.35) and deduce that

$$0 \geq \langle g + Bs_{\circ}, s_{\bullet} - s_{\circ} \rangle + \sigma \|s_{\circ}\| \langle s_{\circ}, s_{\bullet} - s_{\circ} \rangle,$$

which, together with (4.42), gives that

$$\begin{aligned} m(x_{\circ}) - m(x_{\bullet}) &\geq \frac{1}{2} \sigma \|s_{\circ}\| \langle s_{\circ}, s_{\bullet} - s_{\circ} \rangle + \frac{1}{2} \sigma \|s_{\bullet}\| \langle s_{\bullet}, s_{\bullet} - s_{\circ} \rangle - \frac{1}{3} \sigma \|s_{\bullet}\|^3 + \frac{1}{3} \sigma \|s_{\circ}\|^3 \\ &\geq \sigma \left(-\frac{1}{2} \|s_{\circ}\|^2 \|s_{\bullet}\| - \frac{1}{6} \|s_{\circ}\|^3 + \frac{1}{6} \|s_{\bullet}\|^3 - \frac{1}{2} \|s_{\bullet}\|^2 \|s_{\circ}\| \right), \end{aligned} \quad (4.43)$$

where we have used the Cauchy–Schwarz inequality. Taking now (4.32) and (4.39) into account and using the fact that $\kappa \leq 1$, we obtain that

$$m(x) - m(x_{\bullet}) \geq m(x_{\circ}) - m(x_{\bullet}) > \kappa \sigma \left(-\frac{1}{2} \alpha^2 - \frac{1}{6} \alpha^3 + \frac{1}{6} - \frac{1}{2} \alpha \right) \|s_{\bullet}\|^3. \quad (4.44)$$

We now select the value of α for which the lower bounds (4.38) and (4.44) are equal, namely $\alpha_* \approx 0.2418$, the only real positive root of $7\alpha^3 + 3\alpha^2 + 3\alpha = 1$. The desired result now follows from (4.38) and (4.44) with $\kappa_{\text{lm}} \stackrel{\text{def}}{=} \alpha_*^3 \approx 0.0141$. \square

Next, we prove the required model decrease under AS8.

LEMMA 4.6 Suppose that AS8 holds at iteration k . Then, there exists a constant $\kappa_{\text{red}} > 0$ independent of k such that (4.25) holds.

Proof. If $l_k = 1$ then the conclusion immediately follows from Lemma 4.4. Otherwise, (4.31) at $i = 1$ and $x_{k,0} = x_k$ imply that Lemma 4.4 applies with $x_k^+ = x_{k,1}$, giving

$$m_k(x_k) - m_k(x_{k,1}) \geq \frac{1}{6} \sigma_k \|x_{k,1} - x_k\|^3.$$

AS8 further implies that

$$m_k(x_{k,2}) \leq m_k(x_{k,1}), \quad \langle \nabla_s m_k(x_{k,2}), x_{k,2} - x_{k,1} \rangle \leq 0, \quad \text{and} \quad \langle \nabla_s m_k(x_{k,1}), x_{k,2} - x_{k,1} \rangle \leq 0.$$

We may then apply Lemma 4.5 a first time with $x_{\circ} = x_{k,1}$ and $x_{\bullet} = x_{k,2}$ to deduce

$$m_k(x_k) - m_k(x_{k,2}) \geq \frac{1}{6} \kappa_{\text{lm}} \sigma_k \|x_{k,2} - x_k\|^3.$$

If $l_k > 2$ we then apply the same technique $l_k - 1$ times: for $i = 2, \dots, l_k$, we deduce from AS8 that

$$m_k(x_{k,i}) \leq m_k(x_{k,i-1}), \quad \langle \nabla_s m_k(x_{k,i}), x_{k,i} - x_{k,i-1} \rangle \leq 0, \quad \text{and} \quad \langle \nabla_s m_k(x_{k,i-1}), x_{k,i} - x_{k,i-1} \rangle \leq 0,$$

while we obtain by induction that

$$m_k(x_{k,i-1}) \leq m_k(x_k) - \frac{1}{6} \kappa_{\text{lm}}^{i-2} \sigma_k \|x_{k,i-1} - x_k\|^3.$$

This then allows us to apply Lemma 4.5 with $x_{k,\circ} = x_{k,i-1}$ and $x_{k,\bullet} = x_{k,i}$, yielding that

$$m_k(x_k) - m_k(x_{k,i}) \geq \frac{1}{6} \kappa_{\text{lm}}^{i-1} \sigma_k \|x_{k,i} - x_k\|^3.$$

After $\ell_k - 1$ applications of Lemma 4.5 we obtain that

$$m_k(x_k) - m_k(x_{k,\ell_k}) \geq \frac{1}{6} \kappa_{\text{lm}}^{\ell_k-1} \sigma_k \|x_{k,\ell_k} - x_k\|^3. \quad (4.45)$$

Since $x_{k,\ell_k} = x_k^+$ and $s_k = x_k^+ - x_k$, (4.45) is the desired bound (4.25) with $\kappa_{\text{red}} = \frac{1}{6} \kappa_{\text{lm}}^{\bar{\ell}-1}$. \square

4.2.3 Further comments on satisfying AS8. In practice, verifying AS8 need not be too burdensome. Firstly, the computation of x_k^+ could be performed by a sequence of line minimizations, and AS8 then trivially holds provided the number of such minimizations remains uniformly bounded. If the trial step has been determined by another technique one might proceed as follows; see Fig. 1; if we set $x_{k,b}$ to be the global minimizer of the model in the hyperplane orthogonal to the gradient, that is

$$x_{k,b} \stackrel{\text{def}}{=} \underset{\langle g_k, s \rangle = 0}{\text{argmin}} m_k(x_k + s), \quad (4.46)$$

then we may also define $x_{k,a}$ as the intersection of the segment $[x_k, x_{k,b}]$ with the boundary of \mathcal{F} if $x_{k,b} \notin \mathcal{F}$ and as $x_{k,b}$ if $x_{k,b} \in \mathcal{F}$. Similarly we define $x_{k,c}$ as the intersection of the segment $[x_{k,b}, x_k^+]$ with the boundary of \mathcal{F} if $x_{k,b} \notin \mathcal{F}$ and as $x_{k,b}$ if $x_{k,b} \in \mathcal{F}$. We may now verify (4.31) with the set $\{x_k, x_{k,a}, x_{k,c}, x_k^+\}$. If (4.31) fails, then there is a feasible local minimizer of the model along the path

$$\mathcal{P}_k \stackrel{\text{def}}{=} [x_k, x_{k,a}] \cup [x_{k,a}, x_{k,c}] \cup [x_{k,c}, x_k^+] \quad (4.47)$$

(the middle segment being possibly reduced to the point $x_{k,b}$ when it is feasible); further model minimization may then be started from this point—namely from the feasible local minimizer along (4.47)—in order to achieve the termination condition AS7, ignoring the rest of the path and the trial point x_k^+ .

Note that $x_{k,b}$ in (4.46) is the solution of an essentially unconstrained model minimization (in the hyperplane orthogonal to g_k) and thus can be computed at reasonable cost, which makes checking this version of (4.31) acceptable from the computational point of view, especially since $x_{k,b}$ needs to be computed only once even if several x_k^+ must be tested. Clearly, other choices for $x_{k,b}$ are acceptable, as long as a suitable ‘descent path’ \mathcal{P}_k from x_k to x_k^+ can be determined. Note that the purpose of the descent path is to guarantee that the model decrease (4.25) holds, where $s_k = x_k^+ - x_k$; see also (4.45). See Fig. 1 for an illustration of the path \mathcal{P}_k given by (4.47).

4.2.4 The improved complexity bound for COCARC-S We now have all the ingredients needed for the improved function-evaluation complexity result for COCARC-S.

THEOREM 4.7 Suppose that AS1–AS8 and (4.3) hold, and let $\epsilon \in (0, 1]$. Then there exists a constant $\kappa_{\text{df2}} \in (0, 1)$ such that

$$f(x_k) - f(x_{k+1}) \geq \kappa_{\text{df2}} \chi_{k+1}^{3/2} \quad \text{for all } k \in \mathcal{S}. \quad (4.48)$$

Therefore the total number of successful iterations with

$$\min(\chi_k, \chi_{k+1}) > \epsilon \quad (4.49)$$

that occur when applying the COCARC-S algorithm is at most

$$\left\lceil \kappa_{\mathcal{S}_2} \epsilon^{-3/2} \right\rceil \stackrel{\text{def}}{=} I_{\mathcal{S}_2}, \quad (4.50)$$

where $\kappa_{\mathcal{S}_2} \stackrel{\text{def}}{=} (f(x_0) - f_{\text{low}})/\kappa_{\text{df2}}$. Assuming (4.49) holds at $k = 0$, the COCARC-S algorithm takes at most $I_{\mathcal{S}_2} + 1$ successful iterations and evaluations of $\nabla_x f$ (and possibly, of H) until it generates a (first) iterate, say x_{j_2} , such that $\chi_{j_2+1} \leq \epsilon$. Furthermore,

$$j_2 \leq \left\lceil \kappa_{*2} \epsilon^{-3/2} \right\rceil \stackrel{\text{def}}{=} J_2,$$

so that the algorithm takes at most J_2 iterations and objective-function evaluations to terminate with $\chi_{j_2+1} \leq \epsilon$, where

$$\kappa_{*2} \stackrel{\text{def}}{=} \kappa_{\mathcal{S}_2} + (1 + \kappa_{\mathcal{S}_2}) \frac{\log(\sigma_{\max}/\sigma_{\min})}{\log \gamma_1},$$

and where σ_{\max} is defined in (4.15).

Proof. Note that, due to (4.45), the definition of the trial point, namely $x_k^+ = x_k + s_k$, and hence of s_k , does not change even when the path defined by AS8 has more than one segment. Thus, Lemmas 4.3 and 4.6 both apply. Recalling that $f(x_k) = m_k(x_k)$ we obtain from (4.3), (4.25) and (4.16) that

$$f(x_k) - m_k(x_k^+) \geq \frac{1}{6} \sigma_{\min} \kappa_{\text{red}} \kappa_s^3 \chi_{k+1}^{3/2},$$

and thus, from the definition of $k \in \mathcal{S}$, (4.48) follows with $\kappa_{\text{df2}} \stackrel{\text{def}}{=} \frac{1}{6} \eta_1 \sigma_{\min} \kappa_{\text{red}} \kappa_s^3$. Thus, we have

$$f(x_k) - f(x_{k+1}) \geq \kappa_{\text{df2}} \epsilon^{3/2} \quad \text{for all } k \in \mathcal{S} \text{ satisfying (4.48)}. \quad (4.51)$$

Letting $|\mathcal{S}_{\max}|$ denote the number of successful iterations satisfying (4.49), and summing (4.51) over all iterations k from 0 to the last successful iteration satisfying (4.49), it follows from the fact that f does not change on unsuccessful iterations and from AS3a (that $|\mathcal{S}_{\max}| < \infty$ and) that

$$f(x_0) - f_{\text{low}} \geq |\mathcal{S}_{\max}| \kappa_{\text{df2}} \epsilon^{3/2},$$

which gives the bound (4.50). This straightforwardly implies that (4.50) also bounds the number of successful iterations up to j_2 , that conform to (4.1), we denote by $|\mathcal{S}_{j_2}|$. To bound the total number of iterations up to j_2 , let $j = j_2$ in (4.6) and deduce, also from (4.3) and (4.15),

$$|\mathcal{U}_{j_2}| \leq \left\lceil (1 + |\mathcal{S}_{j_2}|) \frac{1}{\log \gamma_1} \log \frac{\sigma_{\max}}{\sigma_{\min}} \right\rceil.$$

This and the bound (4.50) on $|\mathcal{S}_{j_2}|$, as well as $j_2 = |\mathcal{S}_{j_2}| + |\mathcal{U}_{j_2}|$, imply the expression of total bound J_2 on j_2 , recalling also that $\epsilon \leq 1$. \square

This result shows a worst-case complexity result in terms of evaluations of the problem's function that is of the same order as that for the unconstrained case (see [Nesterov & Polyak, 2006](#), or [Cartis *et al.*, 2011b](#)).

Note that global convergence to first-order critical points may be ensured for Algorithm COCARC-S (even without AS5–AS8), if one simply ensures that the steps s_k guarantee a model decrease, which is larger than that obtained at the Cauchy point (as computed by Step 1 of Algorithm COCARC), which means that (2.9) must hold; a very acceptable condition. The convergence analysis presented for Algorithm COCARC thus applies without modification.

Despite not requiring additional evaluations of the problem's nonlinear objective, the subproblem solution and its associated complexity are crucial aspects of an efficient COCARC-S algorithm. In particular, to ensure the better complexity bound of Theorem 4.7, on each iteration k , active-set techniques may be applied starting at x_k to approximately minimize the model $m_k(s)$ in \mathcal{F} along a uniformly bounded number of line segments so as to ensure AS8, until the termination condition (4.13) is satisfied. A minimal and simple such approach is the basic COCARC $_\epsilon$ framework, whose iteration complexity when applied to the model subproblem m_k is addressed in [Cartis *et al.* \(2009, Section 4.3\)](#). In practice, a (much) more efficient active-set technique should be employed; but further investigations into theoretical guarantees of finite termination for such methods are needed, which seem nontrivial to derive in the context of AS7 and AS8 due to the combinatorial aspect of both the (nonconvex) objective and the constraints.

5. Conclusions and perspectives

We have generalized the adaptive cubic regularization method for unconstrained optimization to the case where convex constraints are present. Our method is based on the use of the orthogonal projector onto the feasible domain and is therefore practically limited to situations where applying this projector is computationally inexpensive. This is, for instance, the case if the constraints are simple lower and upper bounds on the variables or if the feasible domain has a special shape such as a sphere, a cylinder or the order simplex (see [Conn *et al.*, 2000, Section 12.1.2](#)). The resulting COCARC algorithm has been proved globally convergent to first-order critical points. This result has capitalized on the natural definition of the first-order criticality measure (3.1), which allows an extension of the unconstrained proof techniques to the constrained case. As a by-product, the Lipschitz continuity of the criticality measure $\chi(x)$ has also been proved for bounded convex feasible sets.

A variant of Algorithm COCARC has then been presented for which a worst-case function-evaluation complexity bound can be shown, which is of the same order as that known for the unconstrained case and better than for steepest descent methods. Remarkably, this algorithm does not rely on global model minimization, but the result obtained is only in terms of the global number of iterations and the problem's function evaluations, leaving aside the complexity of solving the subproblem, even approximately.

The authors are well aware that many issues remain open at this stage, among which the details of an effective step computation, the convergence to second-order points and its associated rate of convergence, and the constraint identification properties, as well as the implications of the new complexity result on optimization problems with equality and inequality constraints. Numerical experience is also necessary to assess the practical potential of both algorithms.

Acknowledgements

The authors are grateful to Ken McKinnon for supplying the main idea for the proof of Lemma 3.3 and to Dimitri Tomanos for carefully reading an early draft of this paper. All three authors are grateful to the Royal Society for its support through the International Joint Project 14265.

Funding

EPSRC grants EP/E053351/1, EP/F005369/1 and EP/G038643/1 and by the European Science Foundation through the OPTPDE program to the third author.

REFERENCES

- CARTIS, C., GOULD, N. I. M. & TOINT, PH. L. (2011a) Adaptive cubic regularisation methods for unconstrained optimization. Part I: motivation, convergence and numerical results. *Math. Program.*, **127**, 245–295.
- CARTIS, C., GOULD, N. I. M. & TOINT, PH. L. (2009) An adaptive cubic regularization algorithm for nonconvex optimization with convex constraints and its function-evaluation complexity. *ERGO Technical Report 09-004*. Edinburgh, UK: School of Mathematics, University of Edinburgh.
- CARTIS, C., GOULD, N. I. M. & TOINT, PH. L. (2011b) Adaptive cubic regularisation methods for unconstrained optimization. Part II: worst-case function- and derivative-evaluation complexity. *Math. Program.*, **130**, 295–319.
- CARTIS, C., GOULD, N. I. M. & TOINT, PH. L. (2010) On the complexity of steepest descent, Newton’s and regularized Newton’s methods for nonconvex unconstrained optimization. *SIAM J. Optimization*, **20**, 2833–2852.
- CONN, A. R., GOULD, N. I. M. & TOINT, PH. L. (1988) Global convergence of a class of trust region algorithms for optimization with simple bounds. *SIAM J. Numer. Anal.*, **25**, 433–460. See also same journal **26**, 764–767.
- CONN, A. R., GOULD, N. I. M. & TOINT, PH. L. (2000) *Trust-Region Methods*. Number 01 in ‘MPS-SIAM Series on Optimization’. Philadelphia, PA: SIAM.
- CONN, A. R., GOULD, N. I. M., SARTENAER, A. & TOINT, PH. L. (1993) Global convergence of a class of trust region algorithms for optimization using inexact projections on convex constraints. *SIAM J. Optimization*, **3**, 164–221.
- GRATTON, S., MOUFFE, M., TOINT, PH. L. & WEBER-MENDONÇA, M. (2008a) A recursive trust-region method in infinity norm for bound-constrained nonlinear optimization. *IMA J. Numer. Anal.*, **28**, 827–861.
- GRATTON, S., SARTENAER, A. & TOINT, PH. L. (2008b) Recursive trust-region methods for multiscale nonlinear optimization. *SIAM J. Optimization*, **19**, 414–444.
- GRIEWANK, A. (1981) The modification of Newton’s method for unconstrained optimization by bounding cubic terms. *Technical Report NA/12*. Cambridge, UK: Department of Applied Mathematics and Theoretical Physics, University of Cambridge.
- HIRIART-URRUTY, J.-B. & LEMARÉCHAL, C. (1993) *Convex Analysis and Minimization Algorithms. Part 1: Fundamentals*. Heidelberg: Springer.
- MANGASARIAN, O. L. & ROSEN, J. B. (1964) Inequalities for stochastic nonlinear programming problems. *Oper. Res.*, **12**, 143–154.
- NESTEROV, YU. (2004) *Introductory Lectures on Convex Optimization*. Applied Optimization. Dordrecht, The Netherlands: Kluwer Academic Publishers.
- NESTEROV, YU. (2006) Cubic regularization of Newton’s method for convex problems with constraints. *Technical Report 2006/9*. Belgium: CORE, UCL, Louvain-la-Neuve.
- NESTEROV, YU. & POLYAK, B. T. (2006) Cubic regularization of Newton method and its global performance. *Math. Program.*, **108**, 177–205.
- ROCKAFELLAR, R. T. (1970) *Convex Analysis*. Princeton, NJ: Princeton University Press.

- SARTENAER, A. (1993) Armijo-type condition for the determination of a generalized Cauchy point in trust region algorithms using exact or inexact projections on convex constraints. *Belg. J. Oper. Res. Stat. Comput. Sci.*, **33**, 61–75.
- VAVASIS, S. A. (1991) *Nonlinear Optimization: Complexity Issues*, International Series of Monographs on Computer Science, vol. 8. Oxford: Oxford University Press.
- VAVASIS, S. A. (1993) Black-box complexity of local minimization. *SIAM J. Optimization*, **3**, 60–80.
- WEISER, M., DEUFLHARD, P. & ERDMANN, B. (2007) Affine conjugate adaptive Newton methods for nonlinear elastomechanics. *Optimization Methods and Software*, **22**, 413–431.