



ELSEVIER

Contents lists available at ScienceDirect

Journal of Complexity

journal homepage: www.elsevier.com/locate/jco



Optimality of orders one to three and beyond: Characterization and evaluation complexity in constrained nonconvex optimization[☆]

C. Cartis^{a,*}, N.I.M. Gould^b, Ph.L. Toint^c

^a Mathematical Institute, Oxford University, Oxford OX2 6GG, United Kingdom

^b Numerical Analysis Group, Rutherford Appleton Laboratory, Chilton OX110 QX, United Kingdom

^c Namur Center for Complex Systems (naXys) and Department of Mathematics, University of Namur, 61, rue de Bruxelles, B-5000 Namur, Belgium

ARTICLE INFO

Article history:

Received 29 December 2017

Received in revised form 21 September 2018

Accepted 24 October 2018

Available online 12 November 2018

Keywords:

Nonlinear optimization

Constrained problems

High-order optimality conditions

Complexity theory

ABSTRACT

Necessary conditions for high-order optimality in smooth nonlinear constrained optimization are explored and their inherent intricacy discussed. A two-phase minimization algorithm is proposed which can achieve approximate first-, second- and third-order criticality and its evaluation complexity is analyzed as a function of the choice (among existing methods) of an inner algorithm for solving subproblems in each of the two phases. The relation between high-order criticality and penalization techniques is finally considered, showing that standard algorithmic approaches will fail if approximate constrained high-order critical points are sought.

© 2018 Elsevier Inc. All rights reserved.

1. Introduction

Analyzing the evaluation complexity of algorithms for solving the nonlinear nonconvex optimization problem has been an active research area over the past few years: we refer the interested reader to [1–7,9,12–17,20–25,28–39,42,45–49,51–54] for contributions in this specific area. The main focus of this thriving domain is to give (sometimes sharp) bounds on the number of evaluations of a minimization problem's functions (objective and constraints, if relevant) and their derivatives that

[☆] Communicated by K. Klaus Meer

* Corresponding author.

E-mail addresses: coralia.cartis@maths.ox.ac.uk (C. Cartis), nick.gould@stfc.ac.uk (N.I.M. Gould), philippe.toint@unamur.be (Ph.L. Toint).

are, in the worst case, necessary for the considered algorithms to find an approximate critical point of a certain order. It is not uncommon that such algorithms involve costly internal computations, provided the number of calls to the problem functions is kept as low as possible.

In nearly all cases, complexity bounds are given for the task of finding ϵ -approximate first- or (more rarely) second-order critical points, typically using first- or second-order Taylor models of the objective function in a suitable globalization framework such as those that use trust regions or regularization. Notable exceptions are [1] where ϵ -approximate third-order critical points of unconstrained problems are sought, [6,7,18–20] where ϵ -approximate first-order critical points are considered using Taylor models of order higher than two for unconstrained, convexly-constrained, least-squares and equality-constrained problems, respectively, and [22] where general ϵ -approximate q th order ($q \geq 1$) critical points of convexly constrained optimization are analyzed using Taylor models of degree q .

Because the present contribution focuses on problems involving a mixture of convex inequality and nonlinear equality constraints, it is useful to set the stage by considering earlier research in this constrained framework. In [15], the worst-case evaluation complexity of finding an ϵ -approximate first-order critical point for smooth nonlinear (possibly nonconvex) optimization problems under convex constraints was examined, using methods involving a second-order Taylor model of the objective function. It was then shown that at most $O(\epsilon^{-3/2})$ evaluations of the objective function and its derivatives are needed to compute such an approximate first-order critical point. This result, identical in order to the best known result for the unconstrained case, assumes that the cost of computing a projection onto the convex feasible set is negligible. It comes however at the price of potentially restrictive technical assumptions (see [15] for details). The analysis of [21] then built on this result by first specializing it to convexly constrained nonlinear least-squares and then using the resulting complexity bound in the context of a two-phase algorithm for a problem class involving general constraints. If ϵ_p and ϵ_D are the primal and dual criticality thresholds, respectively, it was shown that at most $O(\epsilon_p^{-1/2} \epsilon_D^{-3/2})$ evaluations of the problem's functions and their gradients are needed to compute an approximate critical point in that case, where the Karush–Kuhn–Tucker (KKT) conditions are scaled to take the size of the Lagrange multipliers into account. Because of the proof of this result is based on the bound for the convex case, it suffers from the same limitations (not to mention an additional constraint on the relative sizes of ϵ_p and ϵ_D , see [21]). Another more general approach was presented in [46] leading to the same complexity bounds, but at the price of solving a subproblem involving the Jacobian of the original nonlinear constraints. The bounds derived in [28] for a trust-funnel algorithm also consider a scaled KKT condition and are of the same order. The worst-case evaluation complexity of constrained optimization problems was also recently analyzed in [6], allowing for high-order derivatives and models in a framework inspired by that of both [7] and [17,21]. At variance with these latter references, this analysis considers unscaled approximate first-order critical points in the sense that such points satisfy the standard unscaled KKT conditions with accuracy ϵ_p and ϵ_D . None of these papers considers ϵ -approximate second-order points for equality constrained problems, except [9] where first- and second-order optimality were proved for trust-region method defined on manifolds.

The goal of this paper is twofold. The first objective is to fill this gap by deriving complexity bounds for finding ϵ -approximate second- and third-order critical points for the inequality/equality-constrained case. The second is to examine higher-order optimality conditions and to expose the intrinsic difficulties that arise for criticality orders beyond three.

Our presentation is organized as follows. Necessary conditions for higher-order criticality for nonlinear optimization problems involving both convex set constraints and (possibly) nonlinear equality constraints are proposed and discussed in Section 2. A new two-phase algorithm is then introduced in Section 3, whose purpose is to compute ϵ -approximate critical points of orders one and two for such problems, and its evaluation complexity is analyzed in Section 4 as a function of the complexity of an underlying inner algorithm for solving subproblems occurring in each of the two phases. A discussion of the results and some conclusions are finally presented in Section 5.

Basic notation. Our notation is as follows. $y^T x$ denotes the Euclidean inner product of the vectors x and y of \mathbb{R}^n and $\|x\| = (x^T x)^{1/2}$ is the associated Euclidean norm. The cardinality of the set S is denoted by $|S|$. If T_1 and T_2 are tensors, $T_1 \otimes T_2$ is their tensor product and $\|T\|_q$ is the recursively induced Euclidean

(or spectral) norm of the q th order tensor T . If T is a symmetric tensor of order q , the q -kernel of the multilinear q -form

$$T[v]^q \stackrel{\text{def}}{=} T[\underbrace{v, \dots, v}_q \text{ times}]$$

is denoted $\ker^q[T] \stackrel{\text{def}}{=} \{v \in \mathbb{R}^n \mid T[v]^q = 0\}$ (see [10,11]). Note that, in general, $\ker^q[T]$ is a union of cones.¹ If \mathcal{X} is a closed set, \mathcal{X}^0 denotes its interior. The vectors $\{e_i\}_{i=1}^n$ are the coordinate vectors in \mathbb{R}^n . If $\{a_k\}$ and $\{b_k\}$ are two infinite sequences of positive scalars converging to zero, we say that $a_k = o(b_k)$ if and only if $\lim_{k \rightarrow \infty} a_k/b_k = 0$. The normal cone to a general convex set \mathcal{C} at $x \in \mathcal{C}$ is defined by

$$\mathcal{N}_{\mathcal{C}}(x) \stackrel{\text{def}}{=} \{s \in \mathbb{R}^n \mid s^T(z - x) \leq 0 \text{ for all } z \in \mathcal{C}\}$$

and its polar, the tangent cone to \mathcal{C} at x , by

$$\mathcal{T}_{\mathcal{C}}(x) = \mathcal{N}_{\mathcal{C}}^*(x) \stackrel{\text{def}}{=} \{s \in \mathbb{R}^n \mid s^T v \leq 0 \text{ for all } v \in \mathcal{N}_{\mathcal{C}}(x)\}.$$

Note that $\mathcal{C} \subseteq \mathcal{T}_{\mathcal{C}}(x)$ for all $x \in \mathcal{C}$. We also define $P_{\mathcal{C}}[\cdot]$ to be the orthogonal projection onto \mathcal{C} (for an introduction to the relevant properties of convex sets and cones, see [40, Chapter 3] or [50, Part I] for an in-depth treatment.)

2. Necessary optimality conditions for constrained optimization

We consider the smooth constrained problem in the form

$$\min_{x \in \mathcal{F}} f(x) \quad \text{subject to} \quad c(x) = 0 \tag{2.1}$$

where $c : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is sufficiently smooth and f and $\mathcal{F} \subseteq \mathbb{R}^n$ is a non-empty, closed convex set. Note that this formulation covers the problems involving both equality and inequality constraints, the latter being handled using slack variables and the inclusion of the associated simple bounds in the definition of \mathcal{F} .

We start by investigating the necessary optimality conditions for problem (2.1) at x_* expressed in terms of the derivatives of the problem’s function at x_* . Our technique considers paths $x(\alpha)$ of the form

$$x(\alpha) = x_* + \sum_{i=1}^j \alpha^i s_i + o(\alpha^j) \tag{2.2}$$

for $j \in \{1, \dots, q\}$ and where $\alpha > 0$.

Definition: We say that the path $x(\alpha)$ is \mathcal{F} -feasible if $x(\alpha) \in \mathcal{F}$ for all α sufficiently small, and that it is feasible when it is \mathcal{F} -feasible and $c(x(\alpha)) = 0$ for all α sufficiently small.

We define the j th order descriptor set of \mathcal{F} at x by

$$\mathcal{D}_{\mathcal{F}}^j(x) \stackrel{\text{def}}{=} \left\{ (s_1, \dots, s_j) \in \mathbb{R}^{n \times j} \mid x(\alpha) = x + \sum_{i=1}^j \alpha^i s_i + o(\alpha^j) \in \mathcal{F} \right\}, \tag{2.3}$$

which is the set of all (s_1, \dots, s_j) such that a corresponding \mathcal{F} -feasible path $x(\alpha)$ exists. It results from this definition that $\mathcal{D}_{\mathcal{F}}^1(x) = \mathcal{T}_{\mathcal{F}}(x)$, the standard tangent cone to \mathcal{F} at x , and that

$$\mathcal{D}_{\mathcal{F}}^j(x) = \{(s_1, \dots, s_j) \mid (s_1, \dots, s_{j-1}) \in \mathcal{D}_{\mathcal{F}}^{j-1}(x) \text{ and } s_j \in \mathcal{T}_{\mathcal{F}}^{\ell}(x, s_1, \dots, s_{j-1})\} \tag{2.4}$$

¹ The 1-kernels are not only unions of cones but also subspaces. However this is not true for general q -kernels, since both $(0, 1)^T$ and $(1, 0)^T$ belong to the 2-kernel of the non-negative symmetric 2-form $x_1 x_2$ on \mathbb{R}^2 , but their sum does not. $\ker^1[x]$ is the usual orthogonal complement to the vector x , $\ker^2[M]$ is the standard nullspace of the matrix M .

Table 2.2
The sets $\mathcal{P}(j, k)$ for $k \leq j \leq 4$.

$j \downarrow$	$k \rightarrow$			
	1	2	3	4
1	{(1)}			
2	{(2)}	{(1,1)}		
3	{(3)}	{(1,2),(2,1)}	{(1,1,1)}	
4	{(4)}	{(1,3),(2,2),(3,1)}	{(1,1,2),(1,2,1),(2,1,1)}	{(1,1,1,1)}

where, for $j \geq 2$, the ℓ -th order tangent cone to \mathcal{F} at x in the directions s_1, \dots, s_{j-1} is defined by

$$\mathcal{T}_{\mathcal{F}}^j(x, s_1, \dots, s_{j-1}) \stackrel{\text{def}}{=} \{s_\ell \in \mathbb{R}^n \mid x(\alpha) = x + \sum_{k=1}^{j-1} \alpha^k s_k + \alpha^j s_j + o(\alpha^j) \in \mathcal{F}\},$$

generalizing the notion of second-order tangent cone of [8], itself inspired by [26,43]. We also have that

$$s_i \in \mathcal{T}_F(x) \text{ for } i \in \{1, \dots, i_0\}, \tag{2.5}$$

where i_0 is the index of the first nonzero s_i , if any, or $i_0 = j$ otherwise. The necessary optimality conditions for problem (2.1) that we seek involve the index sets $\mathcal{P}(j, k)$ defined, for $k \leq j$, by

$$\mathcal{P}(j, k) \stackrel{\text{def}}{=} \{(\ell_1, \dots, \ell_k) \in \{1, \dots, j\}^k \mid \sum_{i=1}^k \ell_i = j\}. \tag{2.6}$$

For $k \leq j \leq 4$, these are given by Table 2.2.

Theorem 2.1. Suppose that f and each of the $\{c_i\}_{i=1}^m$ are q times continuously differentiable in an open set containing \mathcal{F} , and that x_* is a local minimizer for problem (2.1). Then $c(x_*) = 0$ and, for all feasible paths $x(\alpha)$ of the form (2.2) (if they exist), we have that

$$\sum_{k=1}^j \frac{1}{k!} \left(\sum_{(\ell_1, \dots, \ell_k) \in \mathcal{P}(j, k)} \nabla_x^k f(x_*)[s_{\ell_1}, \dots, s_{\ell_k}] \right) \geq 0 \tag{2.7}$$

for all $j \in \{1, \dots, q\}$ such that

$$\sum_{k=1}^i \frac{1}{k!} \left(\sum_{(\ell_1, \dots, \ell_k) \in \mathcal{P}(i, k)} \nabla_x^k f(x_*)[s_{\ell_1}, \dots, s_{\ell_k}] \right) = 0, \quad (i = 1, \dots, j - 1). \tag{2.8}$$

Moreover, we have that, for all $j \in \{1, \dots, q\}$,

$$\sum_{k=1}^i \frac{1}{k!} \left(\sum_{(\ell_1, \dots, \ell_k) \in \mathcal{P}(i, k)} \nabla_x^k c(x_*)[s_{\ell_1}, \dots, s_{\ell_k}] \right) = 0, \quad (i = 1, \dots, j). \tag{2.9}$$

Proof. Consider a feasible path of the form (2.2). Substituting this relation in the expression $f(x(\alpha)) \geq f(x_*)$ (which must be true for small $\alpha > 0$ if x_* is a local minimizer) and collecting terms of equal degree in α , we obtain that, for sufficiently small α ,

$$0 \leq f(x(\alpha)) - f(x_*) = \sum_{j=1}^q \alpha^j \sum_{k=1}^j \frac{1}{k!} \left(\sum_{(\ell_1, \dots, \ell_k) \in \mathcal{P}(j, k)} \nabla_x^k f(x_*)[s_{\ell_1}, \dots, s_{\ell_k}] \right) + o(\alpha^q) \tag{2.10}$$

where $\mathcal{P}(i, k)$ is defined in (2.6). For this to be true, we need each coefficient of α^j to be non-negative on the zero set of the coefficients $1, \dots, j - 1$ (i.e., satisfying (2.8)), which proves (2.7)–(2.8).

Similarly, substituting (2.2) in the expression $c(x(\alpha)) = 0$ and collecting terms of equal degree in α , we obtain that, for sufficiently small α ,

$$0 = c(x(\alpha)) = \sum_{j=1}^q \alpha^j \sum_{k=1}^j \frac{1}{k!} \left(\sum_{(\ell_1, \dots, \ell_k) \in \mathcal{P}(j,k)} \nabla_x^k c(x_*)[s_{\ell_1}, \dots, s_{\ell_k}] \right) + o(\alpha^q); \tag{2.11}$$

This imposes each coefficient of α^j to be identically zero, yielding (2.9). \square

As is clear from its proof, this theorem may be interpreted as stating necessary optimality conditions for the unidimensional problem of minimizing $f(x(\alpha))$ subject to $c(x(\alpha)) = 0$. The necessary conditions stated by Theorem 2.1 remain very implicit in that they do not directly show the dependence on \mathcal{F} . This dependence is however present since we only consider feasible paths, which must be \mathcal{F} -feasible. Our focus on feasible paths also means that Lagrange multipliers do not appear. Our next steps will bring us closer to standard optimality conditions, but they require a constraint qualification assumption because Theorem 2.1 does not ensure the existence of feasible paths.

AS.0 (Constraint Qualification) Let $x_* \in \mathcal{F}, j \in \{1, \dots, q\}$ and $(s_1, \dots, s_j) \in \mathcal{D}_{\mathcal{F}}^j(x_*)$ be given such that (2.9) holds. Then there exists a feasible path of the form (2.2).

When $q = j = 1$ and \mathcal{F} is described by a set of inequalities, AS.0 is implied by the LICQ first-order constraint qualification.² This can be shown using the implicit function theorem: assuming that the set \mathcal{F} is defined by a set of inequalities, defining $c_{\mathcal{A}}(x)$ to be the subset of equality and inequality constraints active at x , and given s_1 such that $\nabla_x^1 c_{\mathcal{A}}(x_*)[s_1] = 0$, there exists a feasible path with tangent s_1 .

We now introduce the Lagrangian function associated with (2.1)

$$\Lambda(x, y) \stackrel{\text{def}}{=} f(x) + y^T c(x), \tag{2.12}$$

where $y \in \mathbb{R}^m$, and the subspace

$$\mathcal{M}(x) \stackrel{\text{def}}{=} \ker^1[\nabla_x^1 c(x)] \cap \ker^1[\nabla_x f(x)]. \tag{2.13}$$

Theorem 2.2. Suppose that f and each of the $\{c_i\}_{i=1}^m$ are q times continuously differentiable in an open set containing \mathcal{F} , and that x_* is a local minimizer for problem (2.1) at which AS.0 holds. Then we have that $c(x_*) = 0$ and, for all $y_* \in \mathbb{R}^m$ and all $j \in \{1, \dots, q\}$,

$$\sum_{k=1}^j \frac{1}{k!} \left(\sum_{(\ell_1, \dots, \ell_k) \in \mathcal{P}(j,k)} \nabla_x^k \Lambda(x_*, y_*)[s_{\ell_1}, \dots, s_{\ell_k}] \right) \geq 0 \tag{2.14}$$

for all $(s_1, \dots, s_j) \in \mathcal{D}_{\mathcal{F}}^j(x_*)$ and such that

$$\sum_{k=1}^i \frac{1}{k!} \left(\sum_{(\ell_1, \dots, \ell_k) \in \mathcal{P}(i,k)} \nabla_x^k \Lambda(x_*, y_*)[s_{\ell_1}, \dots, s_{\ell_k}] \right) = 0, \quad (i = 1, \dots, j - 1), \tag{2.15}$$

and (2.9) hold. In particular, we have the following.

- If $q = 1$, then

$$\nabla_x^1 \Lambda(x_*, y_*)[s_1] \geq 0 \quad \text{for all } y_* \in \mathbb{R}^m, s_1 \in \mathcal{T}_* \stackrel{\text{def}}{=} \mathcal{T}_{\mathcal{F}}(x_*). \tag{2.16}$$

- If $q = 2$, then

$$\nabla_x^1 \Lambda(x_*, y_*)[s_2] + \frac{1}{2} \nabla_x^2 \Lambda(x_*, y_*)[s_1]^2 \geq 0. \quad \text{for all } y_* \in \mathbb{R}^m, s_1 \in \mathcal{T}_* \cap \mathcal{M}(x_*). \tag{2.17}$$

² Stating that the normals to active constraints are linearly independent.

Proof. AS.0 implies that, if x_* is local minimizer of problem (2.1), then (2.7)–(2.8) in Theorem 2.1 must hold for any of the feasible paths whose existence is implied by any choice of $(s_1, \dots, s_q) \in \mathcal{D}_{\mathcal{F}}^q(x_*)$ subject to (2.9). Adding now (2.7) and (2.8) to the inner product of an arbitrary $y_* \in \mathbb{R}^m$ with (2.9) yields (2.14) and (2.15), respectively. Let us now specialize this result to the case where $q = 1$ (for which conditions (2.15) and (2.9) are void). Observing that $\mathcal{P}(1, 1) = \{(1)\}$ (see Table 2.2), we see that (2.14) implies (2.16). Consider now the case where $q = 2$ and assume that (2.15) and (2.9) hold. The latter implies that

$$s_1 \in \ker^1[\nabla_x^1 c(x_*)],$$

while the former implies that

$$s_1 \in \ker^1[\nabla_x^1 \Lambda(x_*, y_*)]. \tag{2.18}$$

Using these observations and the fact that $s_1 \in \mathcal{T}_*$ because of (2.5), we obtain that the range of s_1 defining feasible paths is further restricted to

$$s_1 \in \mathcal{T}_* \cap \ker^1[\nabla_x^1 c(x_*)] \cap \ker^1[\Lambda(x_*, y_*)] = \mathcal{T}_* \cap \mathcal{M}(x_*).$$

Moreover, using the fact that $\mathcal{P}(2, 1) = \{(2)\}$, $\mathcal{P}(2, 2) = \{(1)\}$ (see Table 2.2) and (2.18), (2.14) then ensures that (2.17) must hold. \square

We note that, as the order j grows, (2.14) and (2.9) for $i = j$ may be interpreted as imposing conditions on s_j (via the derivative tensors of Λ and c at (x_*, y_*) and x_* , respectively), given the directions $\{s_i\}_{i=1}^{j-1}$ satisfying (2.15) and (2.9) for $i \in \{1, \dots, j - 1\}$. Note also that this theorem imposes implicit constraints on y_* in the sense that, for instance, (2.16) ensures that, if s_1 is chosen in \mathcal{T}_* , the inequality on the left of this relation must hold for any y_* , but it does not say that this inequality must hold for any choice of s_1 without restricting the possible y_* . This appears more clearly in the next corollary, which covers some well-known cases and in which the restrictions on s_1 appearing in (2.16) and (2.17) are exchanged for restrictions on y_* .

Corollary 2.3. *Suppose that f and each of the $\{c_i\}_{i=1}^m$ are twice continuously differentiable in an open set containing \mathcal{F} and that x_* is a local minimizer for problem (2.1) at which AS.0 holds. Then we have that $c(x_*) = 0$ and, for some $y_* \in \mathbb{R}^m$,*

$$-\nabla_x^1 \Lambda(x_*, y_*) \in \mathcal{N}_* \stackrel{\text{def}}{=} \mathcal{N}_{\mathcal{F}}(x_*). \tag{2.19}$$

Moreover, if $x_* \in \mathcal{F}^0$, the interior of \mathcal{F} , then $\nabla_x^2 \Lambda(x_*, y_*)$ is positive semi-definite on $\ker^1[\nabla_x^1 c(x_*)]$.

Proof. Using the fact that the normal cone \mathcal{N}_* is the polar of \mathcal{T}_* , we immediately deduce from (2.16) that (2.19) holds. If we also assume that $x_* \in \mathcal{F}^0$, (2.19) unsurprisingly reduces to $\nabla_x^1 \Lambda(x_*, y_*) = 0$, while, for $j = q = 2$, (2.14) gives that $\nabla_x^2 \Lambda(x_*, y_*)$ must be positive semi-definite on the subspace defined by (2.9), that is $\mathcal{M}(x_*) = \ker^1[\nabla_x^1 c(x_*)]$. \square

The conditions stated in Corollary 2.3 for $q = 1$ or 2 are standard (for (2.19), see [27, Theorem 3.2.1, p. 46], for instance, and Fig. 2.1 for an illustration). For more general cases, the complicated conditions (2.14), (2.15) and (2.9) appear not to have been stated before and merit some discussion.

It was observed in [22, Section 3] that the necessary optimality condition for the essentially unconstrained case where $x_* \in \mathcal{F}^0$ (implying $\mathcal{N}_* = \{0\}$) combines more than a single derivative tensor and s_i for orders four and above. If equality constraints are present this situation already appears at order three (and above). Indeed, it can be verified that the necessary conditions (2.14), (2.15) and (2.9) for $q = 3$ and $\mathcal{N}_* = \{0\}$ (and hence $\nabla_x^1 \Lambda(x_*, y_*) = 0$ because of (2.19)) can be written as

$$\nabla_x^2 \Lambda(x_*, y_*)[s_1, s_2] + \frac{1}{6} \nabla_x^3 \Lambda(x_*, y_*)[s_1]^3 = 0 \tag{2.20}$$

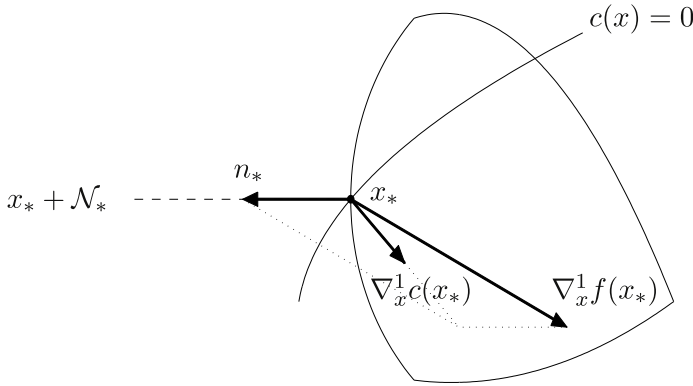


Fig. 2.1. The condition (2.19) with \mathcal{N}_* shown as a dashed half line. Note that $n_* = -\nabla_x^1 \Lambda(x_*, y_*) \neq P_{\mathcal{N}_*}[-\nabla_x^1 f(x_*)]$. Source: adapted from [27].

for all $s_1 \in \mathcal{T}_* \cap \ker^1[\nabla_x^1 \Lambda(x_*, y_*)] \cap \ker^2[\nabla_x^2 \Lambda(x_*, y_*)]$ and

$$\nabla_x^1 c(x_*)[s_2] + \frac{1}{2} \nabla_x^2 c(x_*)[s_1]^2 = 0, \quad \nabla_x^1 c(x_*)[s_3] + \nabla_x^2 c(x_*)[s_1, s_2] + \frac{1}{6} \nabla_x^3 c(x_*)[s_1]^3 = 0. \quad (2.21)$$

These conditions do not require that the second term of the left-hand side of (2.20) vanishes, but instead that it must balance the first term. This is at variance with the unconstrained case, since second-order necessary conditions then ensure that $\nabla_x^2 \Lambda(x_*, y_*)$ is positive semidefinite on \mathbb{R}^n and therefore admits a square root. Thus $\nabla_x^2 \Lambda(x_*, y_*)[s_1, s_2] = [\nabla_x^2 \Lambda(x_*, y_*)^{\frac{1}{2}} s_2]^T [\nabla_x^2 \Lambda(x_*, y_*)^{\frac{1}{2}} s_1] = 0$ since s_1 must belong to $\ker^2[\nabla_x^2 \Lambda(x_*, y_*)]$. However, this argument no longer applies in the constrained case because $\nabla_x^2 \Lambda(x_*, y_*)$ is only positive semidefinite on a strict subspace of \mathbb{R}^n and the square root may fail to exist, as is illustrated by the following example.

Example. Consider the problem

$$\min_{x \in \mathbb{R}^3} x_1 + x_2^2 + x_3^2 - x_3 \quad \text{subject to} \quad c(x) = \begin{pmatrix} -x_1 - x_2^2 + x_1 x_2 + x_3 \\ x_1 + x_2^2 + x_1 x_2 + x_3 \end{pmatrix} = 0,$$

for which the origin is a high-order saddle point.

Comparing the constraints' expression with (2.2) for $q = 3$, we see that

$$s_1 = e_2, \quad s_2 = -e_1 \quad \text{and} \quad s_3 = e_3$$

define a feasible path since then

$$x(\alpha) = \begin{pmatrix} -\alpha^2 \\ \alpha \\ \alpha^3 \end{pmatrix} \quad \text{and} \quad c(x(\alpha)) = \begin{pmatrix} \alpha^2 - \alpha^2 - \alpha^3 + \alpha^3 \\ -\alpha^2 + \alpha^2 - \alpha^3 + \alpha^3 \end{pmatrix} = 0.$$

Now,

$$\nabla_x^1 f(x) = \begin{pmatrix} 1 \\ 2x_2 + 3x_2^2 \\ -1 \end{pmatrix} \quad \nabla_x^2 f(x) = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 2 + 6x_2 & 0 \\ 0 & 0 & 0 \end{pmatrix} \quad \text{and} \quad [\nabla_x^3 f(x)]_{2,2,2} = 6.$$

$$\nabla_x^1 c(x) = \begin{pmatrix} -1 + x_2 & -2x_2 + x_1 & 1 \\ 1 + x_2 & 2x_2 + x_1 & 1 \end{pmatrix}, \quad \nabla_x^2 c_1(x) = \begin{pmatrix} 0 & 1 & 0 \\ 1 & -2 & 0 \\ 0 & 0 & 0 \end{pmatrix}, \quad \nabla_x^3 c_1(x) = 0,$$

$$\nabla_x^2 c_2(x) = \begin{pmatrix} 0 & 1 & 0 \\ 1 & 2 & 0 \\ 0 & 0 & 0 \end{pmatrix} \quad \text{and} \quad \nabla_x^3 c_2(x) = 0.$$

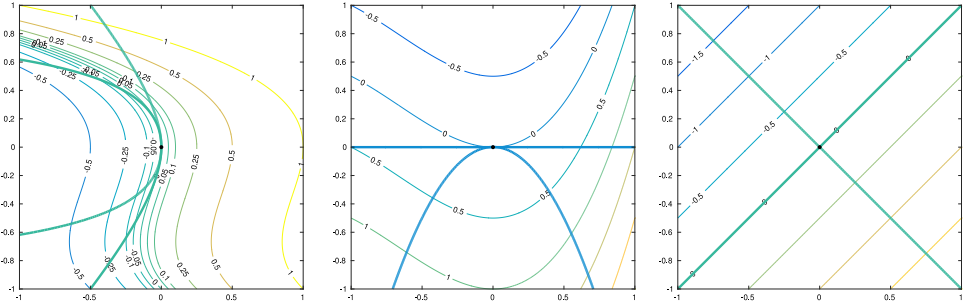


Fig. 2.2. The contour lines of $f(x_1, x_2, 0)$ (left) $f(0, x_2, x_3)$ (center), $f(x_1, 0, x_3)$ (right) and the two constraints intersecting at the origin (thick)

Moreover,

$$\begin{aligned} & \nabla_x^1 c(0)[s_2] + \frac{1}{2} \nabla_x^2 c(0)[s_1]^2 \\ &= - \begin{pmatrix} -1 & 0 & 1 \\ 1 & 0 & 1 \end{pmatrix} e_1 + \frac{1}{2} \left[e_2^T \begin{pmatrix} 0 & 1 & 0 \\ 1 & -2 & 0 \\ 0 & 0 & 0 \end{pmatrix} e_2 \right] e_1 + \frac{1}{2} \left[e_2^T \begin{pmatrix} 0 & 1 & 0 \\ 1 & 2 & 0 \\ 0 & 0 & 0 \end{pmatrix} e_2 \right] e_2 \\ &= 0. \end{aligned}$$

and

$$\begin{aligned} & \nabla_x^1 c(0)[s_3] + \nabla_x^2 c(0)[s_1, s_2] + \frac{1}{6} \nabla_x^3 c(0)[s_1]^3 \\ &= \begin{pmatrix} -1 & 0 & 1 \\ 1 & 0 & 1 \end{pmatrix} e_3 - \left[e_2^T \begin{pmatrix} 0 & 1 & 0 \\ 1 & -2 & 0 \\ 0 & 0 & 0 \end{pmatrix} e_1 \right] e_1 \\ & \quad - \left[e_2^T \begin{pmatrix} 0 & 1 & 0 \\ 1 & 2 & 0 \\ 0 & 0 & 0 \end{pmatrix} e_1 \right] e_2 - \frac{1}{6} 0^T [e_1]^3 \\ &= 0. \end{aligned}$$

Thus (2.21) holds. From the values of $\nabla_x^1 f(0)$ and $\nabla_x^1 c(0)$, we verify that setting $y_0 = (1, 0)^T$ ensures that $\nabla_x^1 \Lambda(0, y_0) = 0$. Hence (2.19) holds as well. Moreover, we have that

$$\ker^1[\nabla_x^1 c(0)] = \ker^1 \left[\begin{pmatrix} -1 & 0 & 1 \\ 1 & 0 & 1 \end{pmatrix} \right] = \text{span} \{e_2\}, \quad \nabla_x^2 \Lambda(0, y_0) = \begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}$$

and the only nonzero component of $\nabla_x^3 \Lambda(0, y_0)$ is its (2,2,2) element which is 6. Thus (2.15) also holds for $i = 2$. In addition, it is easy to check that the third-order necessary condition (2.20) holds with

$$\nabla_x^2 \Lambda(0, y_0)[s_1, s_2] = -1 \quad \text{and} \quad \nabla_x^3 \Lambda(0, y_0)[s_1]^3 = 6.$$

This shows that the term involving $\nabla_x^3 \Lambda(0, y_0)[s_1]^3$ is not the only one occurring in the third-order necessary condition for our example problem, the second derivative term $\nabla_x^2 \Lambda(0, y_0)[s_1, s_2]$ being equally important, as we announced. Fig. 2.2 shows the level lines of the objective function and the constraint manifold in the (x_1, x_2) , (x_2, x_3) and (x_1, x_3) planes, illustrating the interaction of the objective function’s curvature and feasible set. \square

The third order necessary condition therefore must consider both terms in (2.20) and cannot rely only on the third derivative of the Lagrangian along a well-chosen direction or subspace. In general, the q th order necessary conditions will involve (in (2.14)) a mix of other terms than those involving the q th derivative tensor of the Lagrangian applied on vectors s_i ; for $i > 1$, themselves depending on the geometry of the set of feasible arcs. At this stage, for lack of a suitable formal understanding of this geometry, conditions (2.14)–(2.9) remain very difficult to interpret or check.

3. A minimization algorithm

Having analyzed the necessary condition for problem (2.1) and seen that conditions for orders above three are, at this stage, very difficult to verify for general problems, we now describe a two-phase algorithm whose purpose is to find approximate critical points of order one and two (and possibly three as we discuss below). Since the presentation is independent of the order q of the critical points sought, we keep this order general in what follows.

3.1. Inner algorithms for constrained least-squares problems

The new two-phase algorithm relies on an inner algorithm for solving the convexly constrained nonlinear least-squares problem in each of its phases. We therefore start by reviewing the existence and properties of algorithms for solving this subproblem.

Consider first the standard convexly constrained problem

$$\min_{x \in \mathcal{F}} \psi(x) \tag{3.1}$$

where ψ is a smooth function from \mathbb{R}^n to \mathbb{R} and \mathcal{F} is (as in (2.1)) a non-empty closed convex set. An ϵ -approximate q th order critical point for this problem can be defined as a point x such that

$$\phi_{\psi,j}^\Delta(x) \leq \epsilon \Delta^j \text{ for } j = 1, \dots, q \tag{3.2}$$

and some $\Delta \in (0, 1]$, where, for $\mathcal{F}(x) \stackrel{\text{def}}{=} \{d \in \mathbb{R}^n \mid x + d \in \mathcal{F}\}$,

$$\phi_{\psi,j}^\Delta(x) \stackrel{\text{def}}{=} \psi(x) - \underset{\substack{d \in \mathcal{F}(x) \\ \|d\| \leq \Delta}}{\text{globmin}} T_{\psi,j}(x, d), \tag{3.3}$$

is the largest feasible decrease of the j th order Taylor model $T_{\psi,j}(x, s)$ achievable at distance at most Δ from x . Note that $\phi_{\psi,j}^\Delta(x)$ is a continuous function of x and Δ for given \mathcal{F} and f (see [41, Theorem 7]). It is also monotonically increasing in Δ . Also note that the global minimization involved in (3.3) is efficiently solvable for $j = 1$ because it is convex. It is also tractable in the unconstrained case for $j = 2$ since it then reduces to a trust-region subproblem. It may be NP-hard in other cases, which is consistent with the fact that identification of constrained second-order points is NP-hard.

Algorithms for finding ϵ -approximate first-order critical points for problem (3.1), i.e. points satisfying (3.2) for some algorithm-dependent $\Delta \in (0, 1]$ have already been analyzed. Such algorithms, of an essentially theoretical nature, generate a sequence of feasible iterates $\{x_k\}$ with monotonically decreasing objective-function values $\{\psi(x_k)\}$. The method described in [18] proceeds by approximately minimizing models based on the regularized Taylor series of degree p and it can be shown [18, Lemma 2.4]³ that, as long as the stopping criterion (3.2) fails for $q = 1$ and $\Delta = 1$, a sufficient objective-function decrease

$$\psi(x_k) - \psi(x_{k+1}) \geq \kappa_{\text{decr}}^\psi \epsilon^{\frac{p+1}{p}} \tag{3.4}$$

holds for each $k \in \mathcal{S}$, where $\kappa_{\text{decr}}^\psi \in (0, 1)$ is a constant independent of ϵ , and where \mathcal{S} is the set of “successful iterations” at which an effective step is made (i.e. $x_{k+1} \neq x_k$). Moreover, it can also be shown [18, Lemma 2.1] that the set \mathcal{S} cannot be too small in the sense that, for all $k \geq 0$,

$$k \leq \kappa_{\text{uns}}^\psi |\mathcal{S} \cap \{1, \dots, k\}| \tag{3.5}$$

for some constant $\kappa_{\text{uns}}^\psi > 0$. Both $\kappa_{\text{decr}}^\psi$ and κ_{uns}^ψ typically depend on the details of the considered algorithm and of the Lipschitz constant associated with the highest derivative used in the objective-function’s model. Both (3.4) and (3.5) hold under the assumption that $\psi(x)$ is p times continuously differentiable with Lipschitz continuous p th derivative on the “path of iterates” $\cup_{k \geq 0} [x_k, x_{k+1}]$, in that

$$\max_{\xi \in [0,1]} \|\nabla_x^p \psi(x_k + \xi s_k) - \nabla_x^p \psi(x_k)\|_p \leq L_{\psi,p} \|s_k\|, \tag{3.6}$$

³ Observe that $\phi_{\psi,1}^\Delta(x)/\Delta = \chi_{\psi,1}(x)$ as defined in [18, equation (2.4)], irrespective of the value of $\Delta \in (0, 1]$.

for all $\xi \in [0, 1]$, all $k \in \mathcal{S}$ and for some constant $L_{f,p} \geq 0$ independent of x_k and s_k . (Obviously, if the p th derivative of ψ is Lipschitz continuous in an open set containing \mathcal{F} or containing the level set $\{x \in \mathcal{F} \mid \psi(x) \leq \psi(x_0)\}$, then (3.6) holds.)

The algorithm described in [22] is of trust-region type with non-increasing radius. It approximately minimizes a q th degree Taylor inside such a region, Lemma 4.3 in this reference then ensures that, as long as (3.2) fails (for general $q \geq 1$ this time and for Δ being the trust-region radius at iteration k),

$$\psi(x_k) - \psi(x_{k+1}) \geq \kappa_{\text{decr}}^{\psi} \epsilon^{q+1} \tag{3.7}$$

for each $k \in \mathcal{S}$, where we have redefined the constant $\kappa_{\text{decr}}^{\psi}$ to reflect the change in algorithm. In addition, Lemma 4.1 in the same paper also ensures that (3.5) holds for a redefined $\kappa_{\text{uns}}^{\psi}$. Both of these properties again hold if $\psi(x)$ is q times continuously differentiable with Lipschitz continuous q th derivative on the “path of iterates” $\cup_{k \geq 0} [x_k, x_{k+1}]$, in the sense of (3.6) (with p replaced by q).

Summarizing, we see that there exist algorithms for the solution of (3.1) which use truncated Taylor series model of degree q and ensure, under suitable assumptions, both (3.5) and, as long as (3.2) does not hold for some algorithm-dependent non-increasing $\Delta \in (0, 1]$, a lower bound on the objective-function decrease at successful iterations of the form

$$\psi(x_k) - \psi(x_{k+1}) \geq \kappa_{\text{decr}}^{\psi} \epsilon^{\pi} \quad \text{for } k \in \mathcal{S} \tag{3.8}$$

for suitable method-dependent constant $\kappa_{\text{decr}}^{\psi} \in (0, 1)$ and parameter $\pi \geq 1$. (We have that $\pi = (p + 1)/p$ in (3.4) and $\pi = q + 1$ in (3.7).)

Let us now turn to least-squares problems of the form

$$\min_{x \in \mathcal{F}} \psi(x) \stackrel{\text{def}}{=} \frac{1}{2} \|F(x)\|^2, \tag{3.9}$$

(that is problem (3.1) where $\psi(x) = \frac{1}{2} \|F(x)\|_2^2$), where F is a smooth function from \mathbb{R}^n to \mathbb{R}^m . An ϵ -approximate⁴ q th order critical point for this problem can be defined as a point x such that

$$\|F(x)\| \leq \epsilon_p \quad \text{or} \quad \phi_{\psi,j}^{\Delta}(x) \leq \epsilon_D \Delta^j \|F(x)\| \quad \text{for } j = 1, 2 \tag{3.10}$$

and some $\Delta \in (0, 1]$. Note that the second part of (3.10) has the same form as (3.2) with ϵ in the former being replaced by $\epsilon_D \|F(x)\|$ in the latter. It is now easy to verify that, whenever $\|F(x_k)\| \geq \|F(x_{k+1})\|$ and as long as (3.10) fails for x_{k+1} ,

$$\begin{aligned} \|F(x_k)\| (\|F(x_k)\| - \|F(x_{k+1})\|) &\geq \frac{1}{2} (\|F(x_k)\| + \|F(x_{k+1})\|) (\|F(x_k)\| - \|F(x_{k+1})\|) \\ &\geq \frac{1}{2} \|F(x_k)\|^2 - \frac{1}{2} \|F(x_{k+1})\|^2 \\ &= \psi(x_k) - \psi(x_{k+1}) \\ &\geq \kappa_{\text{decr}}^{\psi} [\epsilon_D \|F(x_{k+1})\|]^{\pi}, \end{aligned} \tag{3.11}$$

where we used (3.8) with the form of the second part of (3.10) to derive the last inequality. We will use this last formulation of the guaranteed decrease for least-squares problems as a key piece of our evaluation complexity analysis, together with (3.5) which is needed because the algorithms under consideration require one objective-function evaluation per iteration and one evaluation of its derivatives per successful iteration.

3.2. The outer algorithm

The idea of the two-phase framework which we now introduce is to first apply one of the least-squares algorithms discussed above (or any other method with similar guarantees), which we call Algorithm INNER, to the problem

$$\min_{x \in \mathcal{F}} v(x) \stackrel{\text{def}}{=} \frac{1}{2} \|c(x)\|^2. \tag{3.12}$$

⁴ ϵ_p is the primal accuracy for solving problem (3.9) and ϵ_D the dual one.

(of the form (3.9) with $\psi = \nu$) for finding (under suitably adapted assumptions) an approximate feasible point, if possible. If one is found, Algorithm INNER is then applied to approximately solve the problem

$$\min_{x \in \mathcal{F}} \mu(x, t_k) \stackrel{\text{def}}{=} \frac{1}{2} \|r(x, t_k)\|^2 \stackrel{\text{def}}{=} \frac{1}{2} \left\| \begin{pmatrix} c(x) \\ f(x) - t_k \end{pmatrix} \right\|^2 \tag{3.13}$$

(again of the form (3.9) with $\psi = \mu$) for some monotonically decreasing sequence of “targets” t_k ($k = 1, \dots$). The resulting algorithm is described Algorithm 3.1. Observe that the recomputations of $\phi_{\mu,j}(x_{k+1}, t_{k+1})$ ($j \in \{1, \dots, q\}$) in Step 2.(b) do not require re-evaluating $f(x_{k+1})$ or $c(x_{k+1})$ or any of their derivatives.

Algorithm 3.1. OUTER: a two-phase algorithm for constrained optimization

A starting point x_{-1} and a criticality order $q \in \{1, 2, 3\}$ (for both the feasibility phase and the optimization phase) are given, as well as a constant $\delta \in (0, 1)$. The primal and dual tolerances $0 < \epsilon_p < 1$ and $0 < \epsilon_D < 1$ are also given.

Phase 1: Starting from $x_0 = P_{\mathcal{F}}(x_{-1})$, apply Algorithm INNER to minimize $\nu(x) = \frac{1}{2} \|c(x)\|^2$ subject to $x \in \mathcal{F}$ until a point $x_1 \in \mathcal{F}$ and $\Delta_0 \in (0, 1]$ are found such that

$$\|c(x_1)\| < \delta\epsilon_p \quad \text{or} \quad \phi_{\nu,j}^{\Delta_0}(x_1) \leq \epsilon_D \Delta_0^j \|c(x_1)\| \quad (j \in \{1, \dots, q\}). \tag{3.14}$$

If $\|c(x_1)\| > \delta\epsilon_p$, terminate with $x_\epsilon = x_1$.

Phase 2:

1. Set $t_1 = f(x_1) - \sqrt{\epsilon_p^2 - \|c(x_1)\|^2}$.
2. For $k = 1, 2, \dots$, do:

- (a) Starting from x_k , apply Algorithm INNER to minimize $\mu(x, t_k)$ as a function of $x \in \mathcal{F}$ until an iterate $x_{k+1} \in \mathcal{F}$ and $\Delta_k \in (0, \Delta_{k-1}]$ are found such that

$$\|r(x_{k+1}, t_k)\| < \delta\epsilon_p \quad \text{or} \quad f(x_{k+1}) < t_k \tag{3.15}$$

$$\text{or} \quad \phi_{\mu,j}^{\Delta_k}(x_{k+1}, t_k) \leq \epsilon_D \Delta_k^j \|r(x_{k+1}, t_k)\| \quad (j \in \{1, \dots, q\}).$$

- (b) i. If $\|r(x_{k+1}, t_k)\| < \delta\epsilon_p$, define t_{k+1} according to

$$t_{k+1} = f(x_{k+1}) - \sqrt{\epsilon_p^2 - \|c(x_{k+1})\|^2}. \tag{3.16}$$

and terminate with $(x_\epsilon, t_\epsilon) = (x_{k+1}, t_{k+1})$ if

$$\phi_{\mu,j}^{\Delta_k}(x_{k+1}, t_{k+1}) \leq \epsilon_D \Delta_k^j \|r(x_{k+1}, t_{k+1})\| \quad \text{for } j \in \{1, \dots, q\}. \tag{3.17}$$

- ii. If $\|r(x_{k+1}, t_k)\| \geq \delta\epsilon_p$ and $f(x_{k+1}) < t_k$, define t_{k+1} according to

$$t_{k+1} = 2f(x_{k+1}) - t_k \tag{3.18}$$

and terminate with $(x_\epsilon, t_\epsilon) = (x_{k+1}, t_{k+1})$ if (3.17) holds.

- iii. If $\|r(x_{k+1}, t_k)\| \geq \delta\epsilon_p$ and $f(x_{k+1}) \geq t_k$, terminate with $(x_\epsilon, t_\epsilon) = (x_{k+1}, t_k)$

Algorithm OUTER is again of theoretical interest, and it is unlikely that it would perform well in practice. We now derive some of its useful properties. For this purpose, we partition the Phase 2 outer iterations (before that where termination occurs) into two subsets whose indexes are given by

$$\mathcal{K}_+ \stackrel{\text{def}}{=} \{k \geq 0 \mid \|r(x_{k+1}, t_k)\| < \delta\epsilon_p \quad \text{and} \quad (3.16) \text{ is applied} \} \tag{3.19}$$

and

$$\mathcal{K}_- \stackrel{\text{def}}{=} \{k \geq 0 \mid \|r(x_{k+1}, t_k)\| \geq \delta\epsilon_p \quad \text{and} \quad (3.18) \text{ is applied} \} \tag{3.20}$$

The partition (3.19)–(3.20) allows us to prove then following technical results.

Lemma 3.1. *The sequence $\{t_k\}$ is monotonically decreasing. Moreover, in every Phase 2 iteration of Algorithm OUTER of index $k \geq 1$, we have that*

$$f(x_k) - t_k \geq 0, \tag{3.21}$$

$$\|r(x_{k+1}, t_{k+1})\| = \epsilon_p \text{ for } k \in \mathcal{K}_+, \tag{3.22}$$

$$\|r(x_{k+1}, t_{k+1})\| = \|r(x_{k+1}, t_k)\| \leq \epsilon_p \text{ for } k \in \mathcal{K}_-, \tag{3.23}$$

$$\|c(x_k)\| \leq \epsilon_p \text{ and } f(x_k) - t_k \leq \epsilon_p, \tag{3.24}$$

$$t_k - t_{k+1} \geq (1 - \delta)\epsilon_p \text{ for } k \in \mathcal{K}_+. \tag{3.25}$$

Finally, at termination of Algorithm OUTER,

$$\|r(x_\epsilon, t_\epsilon)\| \geq \delta\epsilon_p, \quad f(x_\epsilon) \geq t_\epsilon \tag{3.26}$$

$$\text{and } \phi_{\mu,j}^{\Delta_k}(x_\epsilon, t_\epsilon) \leq \epsilon_D \Delta_k^q \|r(x_\epsilon, t_\epsilon)\| \text{ for } j \in \{1, \dots, q\}.$$

Proof. The inequality (3.21) follows from (3.16) for $k - 1 \in \mathcal{K}_+$ and from (3.18) for $k - 1 \in \mathcal{K}_-$. (3.22) is also deduced from (3.16) while (3.18) implies the equality in (3.23), the inequality in that statement resulting from the monotonically decreasing nature of $\|r(x, t_k)\|$ during inner iterations in Step 2.(a) of Algorithm OUTER. The inequalities (3.24) then follow from (3.21), (3.22) and (3.23). We now prove (3.25), which only occurs when $\|r(x_{k+1}, t_k)\| \leq \delta\epsilon_p$, that is when

$$(f(x_{k+1}) - t_k)^2 + \|c(x_{k+1})\|^2 \leq \delta^2 \epsilon_p^2. \tag{3.27}$$

From (3.16), we then have that

$$t_k - t_{k+1} = -(f(x_{k+1}) - t_k) + \sqrt{\|r(x_k, t_k)\|^2 - \|c(x_{k+1})\|^2}. \tag{3.28}$$

Now taking into account that the global minimum of the problem

$$\min_{(f,c) \in \mathbb{R}^2} \vartheta(f, c) \stackrel{\text{def}}{=} -f + \sqrt{\epsilon_p^2 - c^2} \text{ subject to } f^2 + c^2 \leq \omega^2,$$

for $\omega \in [0, \epsilon_p]$ is attained at $(f_*, c_*) = (\omega, 0)$ and it is given by $\vartheta(f_*, c_*) = \epsilon_p - \omega$ (see [21, Lemma 5.2]), we obtain from (3.27) and (3.28) (setting $\omega = \delta\epsilon_p$) that

$$t_k - t_{k+1} \geq \epsilon_p - \omega = (1 - \delta)\epsilon_p \text{ for } k \in \mathcal{K}_+$$

for $k \in \mathcal{K}_+$, which is (3.25). Note that, if $k \in \mathcal{K}_-$, then we must have that $t_k > f(x_{k+1})$ and thus (3.18) ensures that $t_{k+1} < t_k$. This observation and (3.25) then allow us to conclude that the sequence $\{t_k\}$ is monotonically decreasing.

In order to prove (3.26), we need to consider, in turn, each of the three possible cases where termination occurs in Step 2.(b). In the first case (i), $\|r(x_{k+1}, t_k)\|$ is small (in the sense that the first inequality in (3.15) holds) and (3.16) is then used, implying that (3.22) holds and that $f(x_{k+1}) > t_{k+1}$. If termination occurs because (3.17) holds, then (3.26) clearly holds at (x_{k+1}, t_{k+1}) . In the second case (ii), the residual $\|r(x_{k+1}, t_k)\|$ is large (the first inequality in (3.15) fails), but $f(x_{k+1}) < t_k$, and t_{k+1} is then defined by (3.18), ensuring that $f(x_{k+1}) > t_{k+1}$ and, because of (3.23), that $\|r(x_{k+1}, t_{k+1})\|$ is also large. As before (3.26) holds at (x_{k+1}, t_{k+1}) if termination occurs because (3.17) is satisfied. The third case (iii) is when $\|r(x_{k+1}, t_k)\|$ is sufficiently large and $f(x_{k+1}) \geq t_k$. But (3.15) then guarantees that $\phi_{\mu,j}^{\Delta_k}(x_{k+1}, t_k) \leq \epsilon_D \Delta_k^j \|r(x_{k+1}, t_k)\|$ for $j \in \{1, \dots, q\}$, and the inequalities (3.26) are again satisfied at (x_{k+1}, t_k) . \square

4. Evaluation complexity

In order to state the smoothness assumptions for problem (2.1), we first define, for some parameter $\beta > 0$, the neighbourhood of the feasible set given by

$$C_\beta = \{x \in \mathcal{F} \mid \|c(x)\| \leq \beta\}.$$

We then assume the following.

AS.1 The feasible set \mathcal{F} is closed, convex and non-empty.

AS.2 The function $v(x)$ is smooth enough to ensure that conditions (3.11) and (3.5) hold for Algorithm INNER applied on problem (3.12).

AS.3 The function $\mu(x, t)$ is smooth enough in x to ensure that conditions (3.11) and (3.5) hold for Algorithm INNER applied on problem (3.13), with constants κ_{decr}^μ and κ_{uns}^μ independent of t .

AS.4 There exist constants $\beta \geq \epsilon_p$ and $f_{\text{low}} \in \mathbb{R}$ such that $f(x) \geq f_{\text{low}}$ for all $x \in C_\beta \stackrel{\text{def}}{=} \{x \in \mathcal{F} \mid \|c(x)\| \leq \beta\}$.

AS.2 and AS.3 remain implicit and depend on the particular inner algorithm used (see Section 3.1). For completeness, we now give conditions on the problem's functions f and $\{c_i\}_{i=1}^m$ which allow the transition between assumptions on f and c and the required ones on the Phase 1 and Phase 2 objective functions v and μ .

Lemma 4.1. *Let $p \geq 1$. Assume that f and $\{c_i\}_{i=1}^m$ are p times continuously differentiable and that their derivatives of order one up to p are uniformly bounded and Lipschitz continuous in an open set containing \mathcal{F} . Let the iterations of Algorithm INNER applied to problem (3.12) be indexed by j . Then (3.6) holds for $\nabla_x^q v(x)$ on every segment $[x_j, x_j + s_j]$ ($j \geq 0$) generated by Algorithm INNER during Phase 1 and any $q \in \{1, \dots, p\}$. The same conclusion holds for $\nabla_x^q \mu(x, t)$ on every segment $[x_j, x_j + s_j]$ ($j \geq 0$) generated by Algorithm INNER during Step 2.(a) of Phase 2 and any $q \in \{1, \dots, p\}$, the Lipschitz constant in this latter case being independent of t .*

Proof. Since

$$\nabla_x^q v(x) = \sum_{i=1}^m \left[\sum_{\ell, j > 0, \ell + j = q} \alpha_{\ell, j} \nabla_x^\ell c_i(x) \otimes \nabla_x^j c_i(x) + c_i(x) \nabla_x^q c_i(x) \right]$$

(where $\{\alpha_{\ell, j}\}$ are suitable non-negative and finite coefficients), condition (3.6) is satisfied on the segment $[x_j, x_j + s_j]$ if (i) the derivatives $\{\nabla_x^{\min\{\ell, j\}} c_i(x)\}_{i=1}^m$ are Lipschitz continuous on $[x_j, x_j + s_j]$, (ii) $\{\nabla_x^{\max\{\ell, j\}} c_i(x)\}_{i=1}^m$ are uniformly bounded on $[x_j, x_j + s_j]$, and (iii) we have that

$$\sum_{i=1}^m \|c_i(x_j + \xi s_j) \nabla_x^q c_i(x_j + \xi s_j) - c_i(x_j) \nabla_x^q c_i(x_j)\|_q \leq L_1 \xi \|s_j\| \tag{4.1}$$

for some constant $L_1 > 0$. The first two of these conditions are ensured by the lemma's assumptions. Moreover,

$$\begin{aligned} & \|c_i(x_j + \xi s_j) \nabla_x^q c_i(x_j + \xi s_j) - c_i(x_j) \nabla_x^q c_i(x_j)\|_q \\ & \leq |c_i(x_j + \xi s_j) - c_i(x_j)| \|\nabla_x^q c_i(x_j + \xi s_j)\|_q \\ & \quad + |c_i(x_j)| \|\nabla_x^q c_i(x_j + \xi s_j) - \nabla_x^q c_i(x_j)\|_q \end{aligned}$$

and the first term on the right-hand side is bounded above by $L^2 \xi \|s_j\|$ and the second by $|c_i(x_j)| L \xi \|s_j\|$. Hence (4.1) holds with

$$L_1 = \sum_{i=1}^m (L^2 + |c_i(x_j)|L) \leq mL^2 + m\|c(x_j)\|L \leq mL^2 + m\|c(x_0)\|L$$

because Algorithm INNER ensures that $\|c(x_j)\| \leq \|c(x_0)\|$ for all $j \geq 0$. As a consequence, the lemma's assumptions guarantee that (3.6) holds with the Lipschitz constant

$$m \left[\left(\max_{i=1, \dots, m} \alpha_i \right) L^2 + L^2 + \|c(x_0)\|L \right].$$

We may now repeat, for $\mu(x, t)$ (with fixed t) the same reasoning as above and obtain that condition (3.6) holds for each segment $[x_j, x_j + s_j]$ generated by Algorithm INNER applied in Step 2.(a) of Phase 2, with Lipschitz constant

$$\begin{aligned} m \left[\left(\max_{i=1, \dots, m} \alpha_i \right) L^2 + L^2 + \|c(x_{j,0})\|L \right] + \left(\max_{i=1, \dots, m} \alpha_i \right) L^2 + L^2 + |f(x_{j,0}) - t_j|L \\ \leq (m + 1) \left[L^2 \left(1 + \max_{i=1, \dots, m} \alpha_i \right) + L \right] \stackrel{\text{def}}{=} L_{\mu,p}, \end{aligned}$$

where we have used (3.22) and $\epsilon_p \leq 1$ to deduce the inequality. Note that this constant is independent of t_j , as requested. \square

As the constants κ_{decr}^μ and κ_{uns}^μ in (3.11) and (3.5) directly depend, for the class of inner algorithms considered, on the Lipschitz constants of the derivatives of μ with respect to x , the independence of these with respect to t ensures that κ_{decr}^μ and κ_{uns}^μ are also independent of t , as requested in AS.3.

We now start the evaluation complexity analysis by examining the complexity of Phase 1 of Algorithm OUTER.

Lemma 4.2. *Suppose that AS.1 and AS.2 hold. Then Phase 1 of Algorithm OUTER terminates with an x_1 such that $\|c(x_1)\| \leq \delta\epsilon_p$ or $\phi_{v,q}^{\Delta k} \leq \epsilon \Delta_k^q$ after at most*

$$\left[\kappa_{\text{CC}}^{\|c\|} \|c(x_0)\| \max \left[\epsilon_p^{-1}, \epsilon_p^{1-\pi} \epsilon_D^{-\pi} \right] \right] + 1$$

evaluations of c and its derivatives, where $\kappa_{\text{CC}}^{\|c\|} \stackrel{\text{def}}{=} 2^{-\pi} \kappa_u [\kappa_{\text{decr}}^v]^{-1} \delta^{1-\pi}$ with κ_{decr}^v being the problem-dependent constant defined in (3.11) for the function $v(x)$ corresponding to (3.12).

Proof. First observe that, as long as Algorithm INNER applied to problem (3.12) has not terminated,

$$\|c(x_\ell)\| \geq \delta\epsilon_p, \tag{4.2}$$

because of the first part of (3.14). Let $\ell \in S_k$ be the index of a successful iteration of Algorithm INNER before termination and suppose first that $\|c(x_{\ell+1})\| \leq \frac{1}{2} \|c(x_\ell)\|$. Then

$$\|c(x_\ell)\| - \|c(x_{\ell+1})\| \geq \frac{1}{2} \|c(x_\ell)\| \geq \frac{1}{2} \delta\epsilon_p \tag{4.3}$$

Suppose now that $\|c(x_{\ell+1})\| > \frac{1}{2} \|c(x_\ell)\|$. As a consequence, we obtain that

$$(\|c(x_\ell)\| - \|c(x_{\ell+1})\|) \|c(x_\ell)\| \geq \kappa_{\text{decr}}^v (\epsilon_D \|c(x_{\ell+1})\|)^\pi$$

where we have also the fact that $\phi_{v,j}^{\Delta k}(x_{\ell+1}) > \epsilon_D \|c(x_{\ell+1})\| \Delta_k^j$ since ℓ occurs before termination, the fact that $\|c(x_\ell)\| \geq \|c(x_{\ell+1})\|$ for $\ell \in S$ and condition (3.11). Hence, using (4.2), we have that

$$\|c(x_\ell)\| - \|c(x_{\ell+1})\| \geq \kappa_{\text{decr}}^v 2^{-\pi} \|c(x_\ell)\|^{\pi-1} \epsilon_D^\pi \geq 2^{-\pi} \kappa_{\text{decr}}^v \delta^{\pi-1} \epsilon_p^{\pi-1} \epsilon_D^\pi.$$

Because of the definition of κ_{decr}^v in (3.11), we thus obtain from this last bound and (4.3) that, for all j ,

$$\|c(x_\ell)\| - \|c(x_{\ell+1})\| \geq \frac{1}{2} \kappa_{\text{decr}}^v \delta^{\pi-1} \min \left[\epsilon_p, \epsilon_p^{\pi-1} \epsilon_D^\pi \right].$$

We then deduce that

$$|S_k| \leq 2[\kappa_{\text{decr}}^v]^{-1} \delta^{-\frac{1}{p}} \|c(x_0)\| \max \left[\epsilon_p^{-1}, \epsilon_p^{1-\pi} \epsilon_D^{-\pi} \right]$$

The desired conclusion then follows by using condition (3.5) and adding one for the final evaluation at termination. \square

Using the results of this lemma allows us to bound the number of outer iterations in \mathcal{K}_+ .

Lemma 4.3. *Suppose that AS.4 holds. Then*

$$|\mathcal{K}_+| \leq \frac{f(x_1) - f_{\text{low}} + 1}{1 - \delta} \epsilon_p^{-1}.$$

Proof. We first note that (3.22) and (3.23) and AS.4 ensure that $x_k \in \mathcal{C}_\beta$ for all $k \geq 0$. The result then immediately follows from AS.4 again and the observation that, from (3.25), t_k decreases monotonically with a decrease of at least $(1 - \delta)\epsilon_p$ for $k \in \mathcal{K}_+$. \square

Consider now x_k for $k \in \mathcal{K}_+$ and denote by $x_{n(k)}$ the next iterate such that $n(k) \in \mathcal{K}_+$ or the algorithm terminates at $n(k)$. Two cases are then possible: either a single pass in Step 2.(a) of Algorithm OUTER is sufficient to obtain $x_{n(k)}$ ($n(k) = k + 1$) or two or more passes are necessary, with iterations $k + 1, \dots, n(k) - 1$ belonging to \mathcal{K}_- . Assume now that the iterations of Algorithm INNER at Step 2.(a) of the outer iteration ℓ are numbered $(\ell, 0), (\ell, 1), \dots, (\ell, e_\ell)$ and note that the mechanism of Algorithm OUTER ensures that iteration (ℓ, e_ℓ) is successful for all ℓ . Now define, for $k \in \mathcal{K}_+$, the index set of all inner iterations necessary to deduce $x_{n(k)}$ from x_k , that is

$$\mathcal{I}_k \stackrel{\text{def}}{=} \{(k, 0), \dots, (k, e_k), \dots, (\ell, 0), \dots, (\ell, e_\ell), \dots, (n(k) - 1, 0), \dots, (n(k) - 1, e_{n(k)-1})\} \quad (4.4)$$

where $k < \ell < n(k) - 1$. Observe that, by the definitions (3.19) and (4.4), the index set of all inner iterations before termination is given by $\cup_{k \in \mathcal{K}_+} \mathcal{I}_k$, and therefore that the number of evaluations of problem’s functions required to terminate in Phase 2 is bounded above by

$$\left| \bigcup_{k \in \mathcal{K}_+} \mathcal{I}_k \right| + 1 \leq \left(\frac{f(x_1) - f_{\text{low}} + 1}{1 - \delta} \epsilon_p^{-1} \times \max_{k \in \mathcal{K}_+} |\mathcal{I}_k| \right) + 1, \quad (4.5)$$

where we added 1 to take the final evaluation into account and where we used Lemma 4.3 to deduce the inequality. We now invoke the complexity properties of Algorithm INNER applied to problem (3.13) to obtain an upper bound on the cardinality of each \mathcal{I}_k .

Lemma 4.4. *Suppose that AS.1–AS.3 hold. Then, for each $k \in \mathcal{K}_+$ before termination,*

$$|\mathcal{I}_k| \leq (1 - \delta) \kappa_{\text{CC}}^\mu \max \left[1, \epsilon_p^{2-\pi} \epsilon_D^{-\pi} \right].$$

where κ_{CC}^μ is independent of ϵ_p and ϵ_D and captures the problem-dependent constants associated with problem (3.13) for all values of t_k generated by the algorithm.

Proof. Observe that (3.23) and the mechanism of this algorithm guarantee the strictly decreasing nature of the sequence $\{\|r(x_\ell, t_\ell)\|\}_{\ell=k}^{n(k)-1}$ and hence of the sequence $\{\|r(x_{\ell,s}, t_\ell)\|\}_{(\ell,s) \in \mathcal{I}_k}$. For each $k \in \mathcal{K}_+$, this reduction starts from the initial value $\|r(x_{k,0}, t_k)\| = \epsilon_p$ and is carried out for all iterations with index in \mathcal{I}_k at worst until it is smaller than $\delta\epsilon_p$ (see the first part of (3.15)) or $\phi_{\mu,j}(x_{\ell,s}) \leq \epsilon_D \Delta_k^j \|r(x_{\ell,s+1}, t_\ell)\|$ for $j \in \{1, \dots, q\}$. We may then invoke (3.13) and (3.11) to deduce that, if $(k, s) \in \mathcal{I}_k$,

$$\left(\|r(x_{k,s}, t_k)\| - \|r(x_{k,s+1}, t_k)\| \right) \|r(x_{k,s}, t_k)\| \geq \kappa_{\text{decr}}^\mu (\epsilon_D \|r(x_{k,s+1}, t_k)\|)^\pi, \quad (4.6)$$

for $0 \leq s < e_k$, while

$$\frac{1}{2} \|r(x_{k,e_k}, t_k)\| - \frac{1}{2} \|r(x_{k+1,0}, t_{k+1})\| = 0.$$

As above, suppose first that $\|r(x_{k,s+1}, t_k)\| \leq \frac{1}{2} \|r(x_{k,s}, t_k)\|$. Then

$$\|r(x_{k,s}, t_k)\| - \|r(x_{k,s+1}, t_k)\| \geq \frac{1}{2} \|r(x_{k,s}, t_k)\| \geq \frac{1}{2} \delta \epsilon_p \quad (4.7)$$

because of the first part of (3.15). If $\|r(x_{k,s+1}, t_k)\| > \frac{1}{2} \|r(x_{k,s}, t_k)\|$ instead, then (4.6) implies that

$$\|r(x_{k,s}, t_k)\| - \|r(x_{k,s+1}, t_k)\| \geq \kappa_{\text{decr}}^\mu 2^{-\pi} \|r(x_{k,s}, t_k)\|^{\pi-1} \epsilon_D^\pi.$$

Combining this bound with (4.7) gives that

$$\|r(x_{k,s}, t_k)\| - \|r(x_{k,s+1}, t_k)\| \geq 2^{-\pi} \kappa_{\text{decr}}^\mu \delta^{\pi-1} \min[\epsilon_p, \epsilon_p^{\pi-1} \epsilon_D^\pi].$$

and therefore, as in Lemma 4.2, that

$$|I_k| \leq 2^\pi [\kappa_{\text{decr}}^\mu]^{-1} \delta^{1-\pi} \left[\frac{\epsilon_p - \delta \epsilon_p}{\min[\epsilon_p, \epsilon_p^{\pi-1} \epsilon_D^\pi]} \right] = 2^\pi (1 - \delta) \delta^{1-\pi} [\kappa_{\text{decr}}^\mu]^{-1} \max[1, \epsilon_p^{2-\pi} \epsilon_D^{-\pi}],$$

and the conclusion follows with $\kappa_{\text{CC}}^\mu \stackrel{\text{def}}{=} 2^\pi \delta^{1-\pi} [\kappa_{\text{decr}}^\mu]^{-1}$. \square

We finally combine the above results in a final theorem stating an evaluation complexity bound for Algorithm OUTER in terms of the measures $\phi_{v,j}^{\Delta_k}(x_\epsilon)$.

Theorem 4.5. *Suppose that AS.1–AS.4 hold. Then, for some constants $\kappa_{\text{CC}}^{\|c\|}$ and κ_{CC}^μ independent of ϵ_p and ϵ_D , Algorithm OUTER applied to problem (2.1) needs at most*

$$\left\lceil \left(\kappa_{\text{CC}}^{\|c\|} \|c(x_0)\| + \kappa_{\text{CC}}^\mu [f(x_1) - f_{\text{low}} + 1] \right) \max[\epsilon_p^{-1}, \epsilon_p^{1-\pi} \epsilon_D^{-\pi}] \right\rceil + 2 \tag{4.8}$$

evaluations of f , c and their derivatives up to order p to compute a point x_ϵ and (possibly) a $t_\epsilon \leq f(x_\epsilon)$ such that, when $t_\epsilon = f(x_\epsilon)$,

$$\|c(x_\epsilon)\| > \delta \epsilon_p, \quad \text{and} \quad \phi_{v,j}^{\Delta_k}(x_\epsilon) \leq \epsilon_D \Delta_k^j \|c(x_\epsilon)\| \quad \text{for } j \in \{1, \dots, q\} \tag{4.9}$$

or, when $t_\epsilon < f(x_\epsilon)$,

$$\|c(x_\epsilon)\| \leq \epsilon_p, \quad \text{and} \quad \phi_{\mu,j}^{\Delta_k}(x_\epsilon, t_\epsilon) \leq \epsilon_D \Delta_k^j \|r(x_\epsilon, t_\epsilon)\| \quad \text{for } j \in \{1, \dots, q\}. \tag{4.10}$$

Proof. If Algorithm OUTER terminates in Phase 1, we immediately obtain that (4.9) holds, and Lemma 4.2 then ensures that the number of evaluations of c and its derivatives cannot exceed

$$\left\lceil \kappa_{\text{CC}}^{\|c\|} \|c(x_0)\| \max[\epsilon_p^{-1}, \epsilon_p^{1-\pi} \epsilon_D^{-\pi}] \right\rceil + 1. \tag{4.11}$$

The conclusions of the theorem therefore hold in this case. Let us now assume that termination does not occur in Phase 1. Then Algorithm OUTER must terminate after a number of evaluations of f and c and their derivatives which is bounded above by the upper bound on the number of evaluations in Phase 1 given by (4.11) plus the bound on the number of evaluations of μ given by (4.5) and Lemma 4.4. Using the inequality $q_i \leq q$ and the facts that $\lfloor a \rfloor + \lfloor b \rfloor \leq \lfloor a + b \rfloor$ for $a, b \geq 0$ and $\lfloor a + i \rfloor = \lfloor a \rfloor + i$ for $a \geq 0$ and $i \in \mathbb{N}$, this yields the combined upper bound

$$\left\lceil \kappa_{\text{CC}}^{\|c\|} \|c(x_0)\| \max[\epsilon_p^{-1}, \epsilon_p^{1-\pi} \epsilon_D^{-\pi}] \right\rceil + \left[(1 - \delta) \kappa_{\text{CC}}^\mu \max[1, \epsilon_p^{2-\pi} \epsilon_D^{-\pi}] \times \left\lceil \frac{f(x_1) - f_{\text{low}} + 1}{1 - \delta} \epsilon_p^{-1} \right\rceil \right] + 2,$$

and (4.8) follows. Remember now that (3.26) holds at termination of Phase 2, and therefore that

$$\epsilon_p \geq \|r(x_\epsilon, t_\epsilon)\| \geq \delta \epsilon_p. \tag{4.12}$$

Moreover, we also obtain from (3.26) that

$$\phi_{\mu,j}^{\Delta_k}(x_\epsilon, t_\epsilon) \leq \epsilon_D \Delta_k^j \|r(x_\epsilon, t_\epsilon)\| \quad \text{for } j \in \{1, \dots, q\}. \tag{4.13}$$

Assume first that $f(x_\epsilon) = t_\epsilon$. Then, using the definition of $r(x, t)$, we deduce that, for $j \in \{1, \dots, q\}$,

$$\phi_{v,j}^{\Delta_k}(x_\epsilon) = \phi_{\mu,j}^{\Delta_k}(x_\epsilon) \leq \epsilon_D \Delta_k^j \|c(x_\epsilon)\|$$

and (4.9) is again satisfied because (4.12) gives that $\|c(x_\epsilon)\| = \|r(x_\epsilon, t_\epsilon)\| \geq \delta \epsilon_p$.

If $f(x_\epsilon) > t_\epsilon$ (the case where $f(x_\epsilon) < t_\epsilon$ is excluded by (3.26)), we see that the inequality $\|c(x_\epsilon)\| \leq \|r(x_\epsilon, t_\epsilon)\| \leq \epsilon_p$, and (4.13) imply (4.10). \square

Note that the bound (4.8) is $O(\epsilon^{-(2\pi-1)})$ whenever $\epsilon_p = \epsilon_D = \epsilon$. Also note that we have used the same algorithm for Phase 1 and Phase 2 of Algorithm OUTER, but we could choose to use different methods of complexity π_v and π_μ , respectively, leading a final bound of the form

$$O\left(\max\left[\epsilon_p^{-1}, \epsilon_p^{1-\pi_v} \epsilon_D^{-\pi_v}\right] + \max\left[\epsilon_p^{-1}, \epsilon_p^{1-\pi_\mu} \epsilon_D^{-\pi_\mu}\right]\right). \tag{4.14}$$

Different criticality order may also be chosen for the two phases, leading to variety of possible complexity outcomes.

It is important to note that the complexity bound given by Theorem 4.5 depends linearly on $f(x_1)$, the value of the objective function at the end of Phase 1. Giving an ϵ -independent upper bound on this quantity is in general impossible, but can be done in some case. A trivial bound can of course be obtained if $f(x)$ is bounded in a neighbourhood of the feasible set, that is $\{x \in \mathcal{F} \mid \|c(x)\| \leq \beta\}$ for some $\beta > 0$. This has the advantage of providing a complexity result which is self-contained (in that it only involves problem-dependent quantities), but it is quite restrictive as it excludes, for instance, problems with equality constraints only ($\mathcal{F} = \mathbb{R}^n$) and coercive objective functions. A bound is also readily obtained if the set \mathcal{F} is itself bounded (for instance when the variables are subject to finite lower and upper bounds) or if one assumes that the iterates generated by Phase 1 remain bounded. This may for example be the case if the set $\{x \in \mathbb{R}^n \mid c(x) = 0\}$ is bounded. For specific choices of the convexly-constrained algorithm applied for Phase 1 of Algorithm OUTER, an ϵ_p -dependent bound can finally be obtained without any further assumption. If Phase 1 is solved using the trust-region based algorithm of [22] and x_1 is produced after k_ϵ iterations of this algorithm, we obtain from the definition of the step that $\|s_k\| \leq \Delta_{\max}$ for all $k \geq 1$. In the same spirit, if the regularization algorithm of [20] is used for Phase 1 and x_1 is produced after k_ϵ iterations of this algorithm, we obtain from the proof of Lemma 2.4 in [20] and the definition of successful iterations that

$$v(x_0) \geq v(x_0) - v(x_1) = \sum_{k \in S_{k_\epsilon}} [v(x_k) - v(x_{k+1})] \geq \frac{\eta\sigma_{\min}}{(p+1)!} \sum_{k \in S_{k_\epsilon}} \|s_k\|^{p+1},$$

giving that

$$\|s_k\| \leq \left(\frac{v(x_0)(p+1)!}{\eta\sigma_{\min}}\right)^{\frac{1}{p+1}}.$$

Hence $\|x_1 - x_0\|$ is itself bounded above by this constant times the (ϵ_p -dependent) number of iterations in Phase 1 given by Lemma 4.2. Using the boundedness of the gradient of $v(x)$ on the path of successful iterates implied by AS.2 then ensures (see the Appendix) the (extremely pessimistic) upper bound

$$f(x_1) = f(x_0) + O\left(\max\left[\epsilon_p^{-1}, \epsilon_p^{1-\pi} \epsilon_D^{-\pi}\right]\right). \tag{4.15}$$

Substituting this bound in (4.8) in effect squares the complexity of obtaining (x_ϵ, t_ϵ) .

Assuming that $f(x_1) - f_{\text{low}}$ can be bounded by a constant independent of ϵ_p and ϵ_D , Table 4.3 gives the evaluation complexity bound for achieving first- and second-order optimality for the problem with additional equality constraints, depending on the choice of underlying algorithm for convexly-constrained optimization. In this table, q is the sought criticality order and p is the degree of the Taylor series being used to model the objective function in the inner algorithm. The table also shows that the use of regularized high-degree models for optimality orders beyond one remains to be explored.

We now consider the link between the necessary conditions derived in Section 2 and the results of Theorem 4.5. For future reference, we start by giving the full expressions of the first four derivatives of $\mu(x, t)$ as a function of x :

$$\nabla_x^1 \mu(x, t) = \sum_{i=1}^m c_i(x) \nabla_x^1 c_i(x) + (f(x) - t) \nabla_x^1 f(x), \tag{4.16}$$

$$\nabla_x^2 \mu(x, t) = \sum_{i=1}^m \left[\nabla_x^1 c_i(x) \otimes \nabla_x^1 c_i(x) + c_i(x) \nabla_x^2 c_i(x) \right] + \nabla_x^1 f(x) \otimes \nabla_x^1 f(x) + (f(x) - t) \nabla_x^2 f(x) \tag{4.17}$$

Table 4.3

Evaluation complexity bounds for Algorithm OUTER as a function of the underlying algorithm for convexly-constrained problems, for ϵ -independent $f(x_1) - f_{low}$ and $\epsilon = \epsilon_p = \epsilon_D$.

q	TR-algo	Regularization		
	$(p = q)$	$p = q$	$p = q + 1$	$p \geq q$
1	$O(\epsilon^{-3})$	$O(\epsilon^{-3})$	$O(\epsilon^{-2})$	$O(\epsilon^{-\frac{p+2}{p}})$
2	$O(\epsilon^{-5})$?	?	?
q	$O(\epsilon^{-(2q+1)})$?	?	?

$$\nabla_x^3 \mu(x, t) = \sum_{i=1}^m \left[3 \nabla_x^2 c_i(x) \otimes \nabla_x^1 c_i(x) + c_i(x) \nabla_x^3 c_i(x) \right] + 3 \nabla_x^2 f(x) \otimes \nabla_x^1 f(x) + (f(x) - t) \nabla_x^3 f(x) \quad (4.18)$$

$$\begin{aligned} \nabla_x^4 \mu(x, t) = \sum_{i=1}^m \left[4 \nabla_x^3 c_i(x) \otimes \nabla_x^1 c_i(x) + 3 \nabla_x^2 c_i(x) \otimes \nabla_x^2 c_i(x) + c_i(x) \nabla_x^4 c_i(x) \right] \\ + 4 \nabla_x^3 f(x) \otimes \nabla_x^1 f(x) + 3 \nabla_x^2 f(x) \otimes \nabla_x^2 f(x) + (f(x) - t) \nabla_x^4 f(x) \end{aligned} \quad (4.19)$$

where \otimes denotes the external product.

We finally establish the consequences of Theorem 4.5 in terms of the functions involved in problem (2.1). Because these results make repeated use of Theorem 3.7 in [22], we first recall this proposition.

Theorem 4.6 ([22, Th. 3.7]). *Suppose that ψ , a general objective function, is q times continuously differentiable and that $\nabla_x^q \psi$ is Lipschitz continuous with constant $L_{\psi,q}$ in an open neighbourhood of a point $x_\epsilon \in \mathcal{F}$ of radius larger than Δ_ϵ . Suppose also that, for some ϵ ,*

$$\phi_{\psi,j}^{\Delta_\epsilon}(x_\epsilon) \leq \epsilon \Delta_\epsilon^j \text{ for } j = 1, \dots, q.$$

Then

$$\psi(x_\epsilon + d) \geq \psi(x_\epsilon) - 2\epsilon \Delta_\epsilon^q \text{ for all } d \in \mathcal{F}(x_\epsilon) \text{ such that } \|d\| \leq \left(\frac{q! \epsilon \Delta_\epsilon^q}{L_{\psi,q}} \right)^{\frac{1}{q+1}}.$$

Theorem 4.7. *Suppose that AS.1–AS.4 hold and that, at (x_ϵ, t_ϵ) and for some $\Delta_\epsilon > 0$, conditions (4.9) hold if $f(x_\epsilon) = t_\epsilon$ or conditions (4.10) hold for $\{1, \dots, q\}$ if $f(x_\epsilon) > t_\epsilon$.*

- (i) *If $f(x_\epsilon) = t_\epsilon$ and, for $j \in \{1, \dots, q\}$, $\nabla_x^j v$ is Lipschitz continuous with constant $L_{v,j}$ in a neighbourhood of x_ϵ of radius larger than Δ_ϵ , then, for each $j \in \{1, \dots, q\}$,*

$$\|c(x_\epsilon)\| > \delta \epsilon_p \text{ and } \|c(x_\epsilon + d)\| \geq \|c(x_\epsilon)\| - 2\epsilon_D \|c(x_\epsilon)\| \Delta_\epsilon^j \quad (4.20)$$

for all $d \in \mathcal{F}(x_\epsilon)$ such that

$$\|d\| \leq \left(\frac{j! \epsilon_D \|c(x_\epsilon)\| \Delta_\epsilon^j}{L_{v,j}} \right)^{\frac{1}{j+1}}.$$

- (ii) *If $f(x_\epsilon) > t_\epsilon$, then, for*

$$y_\epsilon = \frac{c(x_\epsilon)}{f(x_\epsilon) - t_\epsilon}, \quad (4.21)$$

one has that

$$\phi_{\Delta,1}^{\Delta_\epsilon}(x_\epsilon, y_\epsilon) \leq \epsilon_D \Delta_\epsilon \|(1, y_\epsilon^T)\| \text{ and } \widehat{\phi}_{\Delta,j}^{\Delta_\epsilon}(x_\epsilon, y_\epsilon) \leq \epsilon_D \Delta_\epsilon^j \|(1, y_\epsilon^T)\| \quad (j = 2, 3), \quad (4.22)$$

where $\widehat{\phi}_{\Lambda_j}^{\Delta_\epsilon}$ differs from $\phi_{\Lambda_j}^{\Delta_\epsilon}$ in that it uses the feasible set $\mathcal{F}(x_\epsilon) \cap \mathcal{M}(x_\epsilon)$ instead of $\mathcal{F}(x_\epsilon)$. Moreover, if f and c have Lipschitz continuous j th derivatives with constants $L_{f,j}$ and $L_{c,j}$, respectively, then

$$\|c(x_\epsilon)\| \leq \delta\epsilon_P \quad \text{and} \quad f(x_\epsilon + d) \geq f(x_\epsilon) - 2\epsilon_P\|y_\epsilon\| - 2\epsilon_D\Delta_\epsilon^j\|(1, y_\epsilon^T)\| \tag{4.23}$$

for all d such that $d \in \mathcal{M}(x_\epsilon) \cap \mathcal{F}(x_\epsilon)$ whenever $j = 2, 3$, $\|c(x_\epsilon + d)\| \leq \epsilon$, and

$$\|d\| \leq \left(\frac{j! \epsilon_D \Delta_\epsilon^j}{\sqrt{2} \max[L_{f,j}, L_{c,j}]} \right)^{\frac{1}{j+1}}, \tag{4.24}$$

Moreover, the second bound in (4.23) can be simplified to

$$f(x_\epsilon + d) \geq f(x_\epsilon) - 2\epsilon_D\Delta_\epsilon^j\|(1, y_\epsilon^T)\| \tag{4.25}$$

for any d such that $d \in \mathcal{M}(x_\epsilon) \cap \mathcal{F}(x_\epsilon)$ whenever $j = 2, 3$, (4.24) holds, and for which $c(x_\epsilon + d) = 0$ or $c(x_\epsilon + d) = c(x_\epsilon)$.

Proof. Consider first the case where $f(x_\epsilon) = t_\epsilon$ (and thus $\|c(x_\epsilon)\| > \delta\epsilon_P$ because of Theorem 4.5). Note that we only need to consider the case where $\|c(x_\epsilon + d)\| \leq \|c(x_\epsilon)\|$. We have that, for $d \in \mathcal{F}(x_\epsilon)$,

$$\|c(x_\epsilon + d)\| - \|c(x_\epsilon)\| = \frac{\|c(x_\epsilon + d)\|^2 - \|c(x_\epsilon)\|^2}{\|c(x_\epsilon + d)\| + \|c(x_\epsilon)\|} \geq \frac{v(x_\epsilon + d) - v(x_\epsilon)}{\|c(x_\epsilon)\|}$$

and the second part of (4.20) then follows from (4.9) and Theorem 4.6 applied to the function v .

Consider now the case where $f(x_\epsilon) > t_\epsilon$ (and thus $\|c(x_\epsilon)\| \leq \epsilon_P$ because of Theorem 4.5). Focus first on the case where $j = 1$. Theorem 4.5 then ensures that

$$\phi_{\mu,1}^{\Delta_\epsilon}(x_\epsilon, t_\epsilon) \leq \epsilon_D\Delta_\epsilon\|r(x_\epsilon, t_\epsilon)\|.$$

Using now (4.21) and

$$\frac{1}{f(x_\epsilon) - t_\epsilon} \nabla_x^1 \mu(x_\epsilon, t_\epsilon) = J(x_\epsilon)^T \frac{c(x_\epsilon)}{f(x_\epsilon) - t_\epsilon} + \nabla_x^1 f(x_\epsilon) = J(x_\epsilon)^T y_\epsilon + \nabla_x^1 f(x_\epsilon) = \nabla_x^1 \Lambda(x_\epsilon, t_\epsilon). \tag{4.26}$$

one has that (4.22) holds for $j = 1$. Moreover, applying Theorem 4.6, we obtain that

$$\Lambda(x_\epsilon + d, y_\epsilon) \geq \Lambda(x_\epsilon, y_\epsilon) - 2\epsilon_D\Delta_\epsilon\|(1, y_\epsilon^T)\|$$

for all $d \in \mathcal{F}(x_\epsilon)$ such that

$$\|d\| \leq \sqrt{\frac{\|(1, y_\epsilon^T)\| \epsilon \Delta_\epsilon}{L_{f,1} + \|y_\epsilon\| L_{c,1}}}.$$

Using now the fact that, for any $a \geq 0$, $\sqrt{2(1 + a^2)} \geq 1 + a$, we obtain that

$$\|(1, y_\epsilon^T)\| \geq \frac{1 + \|y_\epsilon\|}{\sqrt{2}}. \tag{4.27}$$

Hence we deduce that, for all $d \in \mathcal{F}(x_\epsilon)$ satisfying

$$\|d\| \leq \sqrt{\frac{\epsilon \Delta_\epsilon}{\sqrt{2} \max[L_{f,1}, L_{c,1}]}} \tag{4.28}$$

we have that

$$f(x_\epsilon + d) + y_\epsilon^T c(x_\epsilon + d) \geq f(x_\epsilon) + y_\epsilon^T c(x_\epsilon) - 2\epsilon_D\Delta_\epsilon\|(1, y_\epsilon^T)\| \tag{4.29}$$

and hence, using the Cauchy–Schwarz inequality, that

$$f(x_\epsilon + d) \geq f(x_\epsilon) - \|y_\epsilon\| \|c(x_\epsilon) - c(x_\epsilon + d)\| - 2\epsilon_D\Delta_\epsilon\|(1, y_\epsilon)\|.$$

If one additionally requests that $\|c(x_\epsilon + d)\| \leq \epsilon_p$, then, from the first part of (4.10), $\|c(x_\epsilon) - c(x_\epsilon + d)\| \leq 2\epsilon_p$ and therefore $f(x_\epsilon + d) \geq f(x_\epsilon) - 2\epsilon_p \|y_\epsilon\| - 2\epsilon_D \Delta_\epsilon \|(1, y_\epsilon^T)\|$ for all $d \in \mathcal{F}(x_\epsilon)$ such that (4.28) holds. Also note that, if d exists such that $c(x_\epsilon + d) = 0$, $x_\epsilon + d \in \mathcal{F}$ and (4.28) holds, then (4.29) ensures that

$$f(x_\epsilon + d) \geq f(x_\epsilon) - 2\epsilon_D \Delta_\epsilon \|(1, y_\epsilon^T)\| \tag{4.30}$$

since $y_\epsilon^T c(x_\epsilon) \geq 0$ because $f(x_\epsilon) - t_\epsilon > 0$. Similarly, if d exists such that $c(x_\epsilon + d) = c(x_\epsilon)$, $d \in \mathcal{F}(x_\epsilon)$ and (4.28) holds, then (4.29) ensures that (4.30) also holds.

Now turn to the case where $f(x_\epsilon) > t_\epsilon$ and $j = 2$. Observe now that, because of (4.17) and (2.13),

$$\nabla_x^2 \Lambda(x_\epsilon, y_\epsilon)[d]^2 = \frac{1}{f(x_\epsilon) - t_\epsilon} \nabla_x^2 \mu(x_\epsilon, t_\epsilon)[d]^2 \text{ for all } d \in \mathcal{M}(x_\epsilon). \tag{4.31}$$

Now, $\phi_{\mu,2}^{\Delta_\epsilon}(x_\epsilon) \leq \epsilon_D \Delta_\epsilon^2 \|r(x_\epsilon, t_\epsilon)\|$ implies that

$$\nabla_x^1 \mu(x_\epsilon, t_\epsilon)[d] + \frac{1}{2} \nabla_x^2 \mu(x_\epsilon, t_\epsilon)[d]^2 \geq -\epsilon_D \Delta_\epsilon^2 \|r(x_\epsilon, t_\epsilon)\|$$

for all $d \in \mathcal{M}(x_\epsilon) \cap \mathcal{F}(x_\epsilon)$, and thus, dividing by $f(x_\epsilon) - t_\epsilon > 0$ and using (4.26) and (4.31),

$$\nabla_x^1 \Lambda(x_\epsilon, y_\epsilon)[d] + \frac{1}{2} \nabla_x^2 \Lambda(x_\epsilon, y_\epsilon)[d]^2 \geq -\epsilon_D \Delta_\epsilon^2 \|(1, y_\epsilon)\|$$

for all $d \in \mathcal{M}(x_\epsilon) \cap \mathcal{F}(x_\epsilon)$. This in turn ensures that (4.22) holds for $j = 2$. Applying Theorem 4.6 for the problem defining ϕ , we deduce that

$$\Lambda(x_\epsilon + d, y_\epsilon) \geq \Lambda(x_\epsilon, y_\epsilon) - 2\epsilon_D \Delta_\epsilon^2 \|(1, y_\epsilon)\| \tag{4.32}$$

for all d such that $d \in \mathcal{M}(x_\epsilon) \cap \mathcal{F}(x_\epsilon)$. As a consequence, using (4.27) as above, we have that (4.23) holds for $j = 2$ and all $d \in \mathcal{M}(x_\epsilon) \cap \mathcal{F}(x_\epsilon)$ such that

$$\|d\| \leq \left(\frac{2\epsilon_D \Delta_\epsilon^2}{\sqrt{2} \max[L_{f,2}, L_{c,2}]} \right)^{\frac{1}{3}}. \tag{4.33}$$

Applying the same reasoning as above, we deduce that

$$f(x_\epsilon + d) \geq f(x_\epsilon) - 2\epsilon_p \|y_\epsilon\| - 2\epsilon_D \Delta_\epsilon^2 \|(1, y_\epsilon)\|$$

if one additionally requests that $\|c(x_\epsilon + d)\| \leq \epsilon_p$. We may also, as for $j = 1$, deduce from (4.32) that $f(x_\epsilon + d) \geq f(x_\epsilon) - 2\epsilon_D \Delta_\epsilon^2 \|(1, y_\epsilon)\|$ for any d such that $d \in \mathcal{M}(x_\epsilon) \cap \mathcal{F}(x_\epsilon)$ and (4.33) holds and for which $c(x_\epsilon + d) = 0$ or $c(x_\epsilon + d) = c(x_\epsilon)$.

We finally turn to the case where $f(x_\epsilon) > t_\epsilon$ and $j = 3$. It can be verified that, for $s_1 \in \mathcal{M}(x_\epsilon)$,

$$\begin{aligned} \nabla_x^2 \mu(x_\epsilon, t_\epsilon)[s_1, s_2] &= \nabla_x^1 c(x_\epsilon)[s_1] \cdot \nabla_x^1 c(x_\epsilon)[s_2] + \nabla_x^1 f(x_\epsilon)[s_1] \cdot \nabla_x^1 f(x_\epsilon)[s_2] \\ &\quad + (f(x_\epsilon) - t_\epsilon) \nabla_x^2 \Lambda(x_\epsilon, y_\epsilon)[s_1, s_2] \\ &= (f(x_\epsilon) - t_\epsilon) \nabla_x^2 \Lambda(x_\epsilon, y_\epsilon)[s_1, s_2] \end{aligned} \tag{4.34}$$

and

$$\begin{aligned} \nabla_x^3 \mu(x_\epsilon, t_\epsilon)[s_1]^3 &= 3 \left[\sum_{i=1}^m \nabla_x^2 c_i(x_\epsilon)[s_1]^2 \cdot \nabla_x^1 c_i(x_\epsilon)[s_1] + \nabla_x^2 f(x_\epsilon)[s_1]^2 \cdot \nabla_x^1 f(x_\epsilon)[s_1] \right] \\ &\quad + (f(x_\epsilon) - t_\epsilon) \nabla_x^3 \Lambda(x_\epsilon, y_\epsilon)[s_1]^3 \\ &= (f(x_\epsilon) - t_\epsilon) \nabla_x^3 \Lambda(x_\epsilon, y_\epsilon)[s_1]^3. \end{aligned} \tag{4.35}$$

At termination we have that $\phi_{\mu,3}^{\Delta_\epsilon}(x_\epsilon) \leq \epsilon_D \Delta_\epsilon^3 \|r(x_\epsilon, t_\epsilon)\|$, and thus, for all $d \in \mathcal{F}(x_\epsilon)$,

$$\nabla_x^1 \mu(x_\epsilon, t_\epsilon)[d] + \frac{1}{2} \nabla_x^2 \mu(x_\epsilon, t_\epsilon)[d]^2 + \frac{1}{6} \nabla_x^3 \mu(x_\epsilon, t_\epsilon)[d]^3 \geq -\epsilon_D \Delta_\epsilon^3 \|r(x_\epsilon, t_\epsilon)\|.$$

As for $j = 1$ and 2, and for every $d \in \mathcal{M}(x_\epsilon) \cap \mathcal{F}(x_\epsilon)$, the above relations imply that

$$\nabla_x^1 \Lambda(x_\epsilon, y_\epsilon)[d] + \frac{1}{2} \nabla_x^2 \Lambda(x_\epsilon, y_\epsilon)[d]^2 + \frac{1}{6} \nabla_x^3 \Lambda(x_\epsilon, y_\epsilon)[d]^3 \geq -\epsilon_D \Delta_\epsilon^3 \|(1, y_\epsilon^T)\|,$$

and therefore that (4.22) holds for $j = 3$. Applying Theorem 4.6 again, we now deduce that

$$\Lambda(x_\epsilon + d, y_\epsilon) \geq \Lambda(x_\epsilon, y_\epsilon) - 2\epsilon_D \Delta_\epsilon^2 \|(1, y_\epsilon)\|$$

for all $d \in \mathcal{M}(x_\epsilon) \cap \mathcal{F}(x_\epsilon)$. As for the previous cases, this implies that (4.23) holds for $j = 3$ and, using (4.27) once more, for all $d \in \mathcal{M}(x_\epsilon) \cap \mathcal{F}(x_\epsilon)$ satisfying

$$\|d\| \leq \left(\frac{6\epsilon_D \Delta_\epsilon^3}{\sqrt{2} \max[L_{f,3}, L_{c,3}]} \right)^{\frac{1}{3}}. \tag{4.36}$$

The inequality (4.25) is obtained as for the cases where $j = 1, 2$. \square

We verify that (4.22) for $j = 1$ is the scaled first-order criticality condition considered in [21] (Theorem 4.7 thus subsumes the analysis presented in that reference) and is equivalent to

$$\|P_{\mathcal{T}(x_\epsilon)}[-\nabla_x^1 \Lambda(x_\epsilon, y_\epsilon)]\| \leq \epsilon_D \Delta_\epsilon \|(1, y_\epsilon^T)\|,$$

which corresponds to a scaled version of the first-order criticality condition considered in [6].

4.1. Beyond third-order optimality?

We have now proved that, if an approximate q th order critical point for the convexly constrained problem can be obtained by an inner algorithm at a given evaluation complexity, then the same result holds for the critical points of $\|c(x)\|$ whenever Algorithm OUTER terminates with an infeasible stationary point of the constraint violation (either at Phase 1 or at (4.9)). When Algorithm OUTER terminates with (4.10), we have shown in Theorem 4.7 that similar results hold for criticality of orders one, two and three.

As indicated already, the situation becomes considerably more complicated for higher orders. The first difficulty, which we covered in Section 2, is that the conditions (2.14), (2.15) and (2.9) involve, for higher orders, the geometry of the feasible arcs in a way which is hard to exploit. Moreover, the fact that we could derive, in Theorem 4.7, some lower bounds on the objective function values by exploiting information at orders one up to three is strongly dependent of the observation that, in the suitable subspace,

$$\frac{1}{f(x_\epsilon) - t_\epsilon} \nabla_x^j \mu(x_\epsilon, t_\epsilon) = \nabla_x^j \Lambda(x_\epsilon, y_\epsilon) \quad (j = 1, 2, 3) \tag{4.37}$$

(see (4.26), (4.31), (4.34) and (4.35)), which in turn ensures that minimizing $\mu(x, t)$ with respect to x on the said subspace also results in minimizing $\Lambda(x, y)$ with respect to x on the same subspace.⁵ Is this crucial property maintained for high orders? We now show that the answer to this question is negative for orders four and beyond, due to the ever more distant relationship between $\nabla_x^j \mu(x_\epsilon, t_\epsilon)$ and $\nabla_x^j \Lambda(x_\epsilon, y_\epsilon)$ when j grows, which is apparent when considering the expressions (4.16)–(4.19). Indeed, the terms

$$\begin{aligned} \frac{3}{f(x_\epsilon) - t_\epsilon} \left[\sum_{i=1}^m (\nabla_x^2 c_i(x) \otimes \nabla_x^2 c_i(x)) [d]^4 + (\nabla_x^2 f(x) \otimes \nabla_x^2 f(x)) [d]^4 \right] \\ = \frac{3}{f(x_\epsilon) - t_\epsilon} \left[\sum_{i=1}^m (\nabla_x^2 c_i(x) [d]^2)^2 + (\nabla_x^2 f(x) [d]^2)^2 \right] \end{aligned} \tag{4.38}$$

in (4.19) would only vanish in general if $d \in \ker^2[\nabla_x^2 f(x)] \cap \ker^2[\nabla_x^2 c(x)]$. Although this is formally reminiscent of the definition of $\mathcal{M}(x)$ in (2.13), this crucial inclusion now no longer follows from lower-order conditions.

This is illustrated by what happens on the problem

$$\min_{x_1, x_1} -x_2 - x_1^2 + x_1 x_2 - \frac{1}{2} x_1^4 \quad \text{subject to} \quad \varepsilon + x_2 + x_1^2 - x_1 x_2 = 0 \tag{4.39}$$

⁵ For order three, it is fortunate that terms in (4.34) and (4.35) involving the second derivatives always appear in product with terms involving the first, which is the reason why the minimization subspace at order three is not smaller than that at order two.

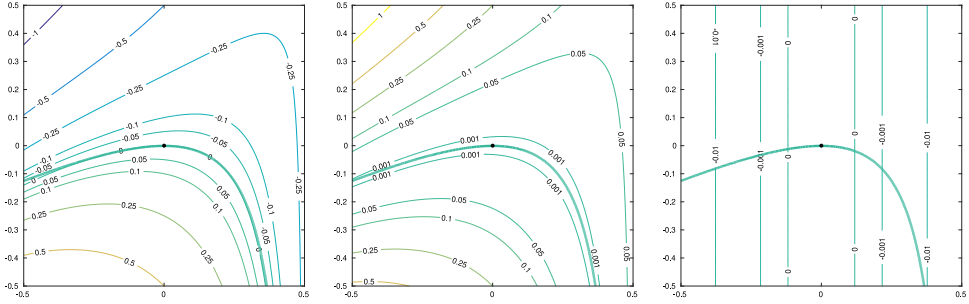


Fig. 4.3. Contour lines for (4.39) (left), $\mu(x, t_\epsilon)$ (center) and $\Lambda(x, y_\epsilon)$ (right) with the constraint shown as a thick curve.

for some $\epsilon \in (0, 1]$. If we consider $x_\epsilon = (0, 0)$ and $t_\epsilon = -\epsilon$ (yielding $y_\epsilon = 1$), then one can verify (see the Appendix) that $\mu(0, t_\epsilon)$ satisfies the necessary conditions for a fourth order minimizer at the origin while the problem itself has a global (fourth order) constrained maximizer. Fig. 4.3 shows the contour lines of the objective function with the constraint set superimposed as a thick curve (left), the contour lines of $\mu(x, t_\epsilon)$ (center) and $\Lambda(x, y_\epsilon)$ (right).

It is worthwhile to note that the above discussion has wider implications. Indeed the first of the problematic terms in (4.38) not only occurs in the function $\mu(x, t)$ used in this paper, but also when applying to problem (2.1) a quadratic, ℓ_1 or ℓ_∞ penalty function, a classical augmented Lagrangian approach, or a sequential quadratic programming method using a merit function depending on such penalty terms. The same difficulty may also occur if more general penalizations of the type $p(v(x))$ (for some increasing smooth function p from \mathbb{R}^+ to \mathbb{R}^+) are employed. Indeed, consider the derivatives of $p(v(x))$. One verifies that

$$\begin{aligned} \nabla_x^4 p(v(x)) &= p''''(v(x))[\nabla_x^1 v(x)]^{4\otimes} + 6p'''(v(x))\nabla_x^2 v(x) \otimes [\nabla_x^1 v(x)]^{2\otimes} \\ &\quad + 4p''(v(x))\nabla_x^3 v(x) \otimes \nabla_x^1 v(x) + 3p'(v(x))[\nabla_x^2 v(x)]^{2\otimes} \\ &\quad + p'(v(x))\nabla_x^4 v(x) \end{aligned}$$

whose last term, together with

$$\nabla_x^4 v(x) = \sum_{i=1}^m \left[4 \nabla_x^3 c_i(x) \otimes \nabla_x^1 c_i(x) + 3 [\nabla_x^2 c_i(x)]^{2\otimes} + c_i(x) \nabla_x^4 c_i(x) \right],$$

indicates that the troublesome terms involving $[\nabla_x^2 c_i(x)]^{2\otimes}$ do not vanish unless $p'(v(x))$ also vanishes with $v(x)$.

None of the linear or quadratic penalization approaches can therefore be expected to reliably produce critical points of orders four or more. Innovative techniques are thus needed if one is interested to compute high-order critical points of (2.1) of higher order. One possible research direction is to follow the propositions formulated in [20] and to exploit penalization terms of order higher than two in the definitions of v and μ , for which an improved evaluation complexity bound is already available for the subproblem solution.

5. Conclusions and discussion

We have formulated and analyzed, in Section 2, the necessary conditions for high-order optimality in nonlinear optimization problems involving both convex set constraints and nonlinear equalities. We have also discussed the difficulties inherent to their form for third-order critical points and higher.

We then have shown in Sections 3 and 4 that the evaluation complexity of finding an approximate q -th-order scaled critical point ($q = 1, 2, 3$) for a large class of smooth nonlinear optimization problem involving both equality and inequality constraints is at most $O(\epsilon_p^{1-\pi} \epsilon_D^{-\pi})$ evaluations of the objective

function, constraints and their derivatives, where ϵ^π is the order of the guaranteed objective function decrease during the successful iterations of an underlying inner algorithm for convexly constrained least-squares problems. We refer here to an “approximate scaled critical point” in that such a point is required to satisfy (4.9) or (4.10), where the accuracy is scaled by the size of the constraint violation or that of the Lagrange multipliers. In particular, the above results provide the first evaluation complexity bound for second- and third-order criticality in the case involving general inequality and equality constraints.

This result also corrects an unfortunate error⁶ in the first-order analysis of [21], that allows a vector of Lagrange multipliers whose sign is arbitrary (in line with a purely first-order setting where minimization and maximization are not distinguished). The present analysis now yields the multiplier with the sign associated with minimization.

Interestingly, an $O(\epsilon_p \epsilon_D^{-(p+1)/p} \min[\epsilon_D, \epsilon_p]^{-(p+1)/p})$ evaluation complexity bound was also proved by Birgin, Gardenghi, Martínez, Santos and Toint in [6] for first-order *unscaled*, standard KKT conditions and in the least expensive of three cases depending on the degree of degeneracy identifiable by the algorithm.⁷ Even if the bounds for the scaled and unscaled cases coincide in order when $\epsilon_p \leq \epsilon_D$, comparing the two results for first-order critical points is not straightforward. On one hand the scaled conditions take into account the possibly different scaling of the objective function and constraints. On the other hand the same scaled conditions may result in earlier termination with (4.10) if the Lagrange multipliers are very large, as (4.10) is then consistent with the weaker requirement of finding a John’s point. But the framework discussed in the present paper also differs from that of [6] in additional significant ways. The first is that second-order critical points are now covered in the analysis. If we now restrict the scope to first-order, the present paper provides a potentially stronger version of the termination of the algorithm at infeasible points (in Phase 1): indeed the second part of (4.9) can be interpreted as requiring that the size of the feasible gradient of $\|c(x)\|$ is below ϵ_D , while [6] considers the gradient of $\|c(x)\|^2$ instead. The second is that, if termination occurs in Phase 2 for an x_ϵ such that $\phi_{v,1}^{\Delta_k}(x_\epsilon)$ is of order $\epsilon_p \epsilon_D \Delta_k$ (thereby covering the case where $f(x_\epsilon) = t_k$ discussed in Theorem 4.5) and $x_\epsilon \in \mathcal{F}^0$, then $\|P_{\mathcal{T}_*}[-\nabla_x^1 v(x_*)]\| = \|\nabla_x^1 v(x_*)\|$ is of the same order and Birgin et al. show that, in this case, the Łojaciewicz inequality [44] must fail for c in the limit for ϵ_p and ϵ_D tending to zero (see [6] for details). This observation is interesting because smooth functions satisfy the Łojaciewicz inequality under relatively weak conditions, implying that termination in these circumstances is unlikely. The same information is also obtained in [6], albeit at the price of worsening the evaluation complexity bound mentioned above by an order of magnitude in ϵ_D . We also note that the approach of [6] requires the minimization, at each iteration, of a residual whose second derivatives are discontinuous, while all functions used in the present paper are p times continuously differentiable. A final difference between the two approaches is obviously our introduction of $\phi_{\mu,j}^{\Delta_k}$ in the expression of the criticality condition in Theorem 4.5 for taking the inequality constraints into account.

Will regularization-based methods provide better evaluation complexity bounds when using polynomial models of higher order? Can the limitations of penalty approaches for finding high-order solutions for equality constrained problems be circumvented? These and many other questions remain open at this stage.

Acknowledgments

The authors would like to thank Oliver Stein for suggesting reference [41]. The work of the second author was supported by EPSRC, United Kingdom grant EP/M025179/1. The third author acknowledges the support provided by the Belgian Fund for Scientific Research (FNRS), the Leverhulme Trust (UK), Balliol College (Oxford, UK), the Department of Applied Mathematics of the Hong Kong Polytechnic University, ENSEIHT (Toulouse, France) and INDAM (Florence, Italy). Thanks are also due to a thoughtful referee whose patience and perceptive comments have helped to significantly improve the manuscript.

⁶ The second equality in the first equation of Lemma 3.4 in [21] only holds if one is ready to flip the gradient’s sign if necessary.

⁷ This result also assumes boundedness of $f(x_1)$.

Appendix

Details of the derivation of (4.15)

For the trust-region algorithm,

$$f(x_1^*) \leq f(x_1) + \left(\max_{\xi \in \cup_{j \in \mathcal{S}} [x_j, x_{j+1}]} \|\nabla_x^1 v(\xi)\| \right) \Delta_{\max} \left(\left[\kappa_{CC}^{\|c\|} \|c(x_1)\| \epsilon_p^{-q} \epsilon_D^{-(q+1)} \right] + 1 \right).$$

For the regularization algorithm,

$$f(x_1) \leq f(x_0) + \left(\frac{v(x_0)(p+1)!}{\eta \sigma_{\min}} \right)^{\frac{1}{p+1}} \times \max_{\xi \in \cup_{j \in \mathcal{S}} [x_j, x_{j+1}]} \|\nabla_x^1 v(\xi)\| \left\{ \left[\kappa_{CC}^{\|c\|} \|c(x_0)\| \max \left[\epsilon_p^{-1}, \epsilon_p^{-\frac{1}{p}}, \epsilon_D^{-\frac{p+1}{p+1-q}} \right] \right] + 1 \right\}.$$

Details for the example (4.39)

We prove the validity of the statement made after the definition of problem (4.39), namely that $\mu(0, t_\epsilon)$ satisfies the necessary conditions for a fourth order minimizer at the origin while the problem itself has a global (fourth order) constrained maximizer.

Let $T_3(x) = x_2 + x_1^2 - 2x_1x_2$ and define, for some $\epsilon \in (0, 1]$,

$$f(x) = -T_3(x) - \frac{1}{2}x_1^4 \quad \text{and} \quad c(x) = \epsilon + T_3(x). \tag{A.1}$$

and thus, for a given multiplier y ,

$$\Lambda(x, y) = -T_3(x) - \frac{1}{2}x_1^4 + y[\epsilon + T_3(x)] \tag{A.2}$$

We have that

$$\nabla_x^1 T_3(x) = \begin{pmatrix} 2x_1 - 2x_2 \\ 1 - 2x_1 \end{pmatrix}, \quad \nabla_x^2 T_3(x) = \begin{pmatrix} 2 & -2 \\ -2 & 0 \end{pmatrix} \quad \text{and} \quad \nabla_x^3 T_3(x) = 0. \tag{A.3}$$

Thus, at the origin and for $t_\epsilon = -\epsilon$

$$c(0) = \epsilon = f(0) - t_\epsilon \quad \text{and} \quad \nabla_x^j c(0) = \nabla_x^j T_3(0) = -\nabla_x^j f(0) \quad \text{for } j = 1, 2, 3. \tag{A.4}$$

As a consequence, the choice $y = 1$, (A.2) and (A.3) ensure that $\Lambda(x, 1) = \epsilon - \frac{1}{2}x_1^4$ as well as

$$\nabla_x^1 \Lambda(0, 1) = 0, \quad \nabla_x^2 \Lambda(0, 1) = 0, \quad \nabla_x^3 \Lambda(0, 1) = 0, \quad \nabla_x^4 \Lambda(0, 1) = \nabla_x^4 f(0) = -12e_1^{\otimes 4}. \tag{A.5}$$

Using (4.16)–(4.19) and (A.4), we also have that, for $t = -\epsilon$,

$$\nabla_x^1 \mu(0, t_\epsilon) = (c(0) - f(0) + t_\epsilon) \nabla_x^1 T_3(0) = (\epsilon - 0 - \epsilon)e_2 = 0, \tag{A.6}$$

$$\nabla_x^2 \mu(0, t_\epsilon) = 2 \nabla_x^1 T_3(0) \otimes \nabla_x^1 T_3(0) = 2e_2 e_2^T, \quad \nabla_x^3 \mu(0, t_\epsilon) = 6 \nabla_x^2 T_3(0) \otimes \nabla_x^1 T_3(0) = 0 \tag{A.7}$$

and, using the last equation in (A.5),

$$\begin{aligned} \nabla_x^4 \mu(0, t_\epsilon) &= 6 \nabla_x^2 T_3(0) \otimes \nabla_x^2 T_3(0) + c(0) \nabla_x^4 c(0) + (f(0) - t_\epsilon) \nabla_x^4 f(0) \\ &= 12 \left[\begin{pmatrix} 1 & -1 \\ -1 & 0 \end{pmatrix}^{\otimes 2} - \epsilon e_1^{\otimes 4} \right]. \end{aligned} \tag{A.8}$$

(Notice the contribution of the first term in the bracketed expression, potentially dwarfing that of the second for sufficiently small ϵ .)

Let us attempt to verify (2.19)–(2.9) with $q = 4$ for the problem of minimizing $\mu(x, t_\epsilon)$ with $s_1 \in \ker^2[\nabla^2\mu] = \text{span}\{e_1\}$. We have that (2.19) holds because of (A.6). We also obtain, from (A.6)–(A.8), that, for $s_1 = \tau e_1$ for some $\tau \in \mathbb{R}$ and for any choice of $s_2, s_3, s_4 \in \mathbb{R}^n$,

$$\begin{aligned} \nabla_x^2 \mu(0, t_\epsilon)[s_2] + \frac{1}{2} \nabla_x^2 \mu(0, t_\epsilon)[s_1]^2 &= 0^T s_2 + \tau e_2 e_2^T e_1 = 0, \\ \nabla_x^3 \mu(0, t_\epsilon)[s_3] + \tau \nabla_x^2 \mu(0, t_\epsilon)[s_1, s_2] + \frac{\tau^3}{6} \nabla_x^3 \mu(0, t_\epsilon)[s_1]^3 &= 0^T s_3 + \tau s_2^T e_2 e_2^T e_1 + \frac{\tau^3}{6} 0[s_1]^3 = 0 \end{aligned}$$

and

$$\begin{aligned} \nabla_x^4 \mu(0, t_\epsilon)[s_4] + \nabla_x^2 \mu(0, t_\epsilon)[s_1, s_3] + \frac{1}{2} \nabla_x^2 \mu(0, t_\epsilon)[s_2]^2 \\ + \frac{1}{2} \nabla_x^3 \mu(0, t_\epsilon)[s_1, s_1, s_2] + \frac{1}{24} \nabla_x^4 \mu(0, t_\epsilon)[s_1]^4 \\ = 0^T s_4 + \tau s_3^T e_2 e_2^T e_1 + \frac{1}{2} (e_2^T s_2)^2 + \frac{\tau^2}{2} 0^T [e_1, e_1, s_2] \\ + \frac{12}{24} \left[\left\| e_1^T \begin{pmatrix} 1 & -1 \\ -1 & 0 \end{pmatrix} e_1 \right\|^2 - \varepsilon \right] \tau^4 (e_1^T e_1)^4 \\ = \frac{1}{2} (e_2^T s_2)^2 + \frac{1}{2} (1 - \varepsilon) \tau^4 (e_1^T e_1)^4. \end{aligned} \tag{A.9}$$

The choice of s_2, s_3 and s_4 is however constrained by (2.9) for $i = 1, 2, 3$, in that those vector must also satisfy the equations

$$\begin{aligned} \nabla_x^1 c(0)[s_1] &= 0 = \tau e_2^T e_1, \\ \nabla_x^1 c(0)[s_2] + \frac{1}{2} \nabla_x^2 c(0)[s_1]^2 &= 0 = e_2^T s_2 + \tau^2 e_1^T \begin{pmatrix} 1 & -1 \\ -1 & 0 \end{pmatrix} e_1, \\ \nabla_x^1 c(0)[s_3] + \nabla_x^2 c(0)[s_1, s_2] + \frac{1}{2} \nabla_x^3 c(0)[s_1]^3 &= 0 = e_2^T s_3 + 2\tau(e_1 - e_2)^T s_2 + \tau^3 0^T [e_1]^3. \end{aligned}$$

and

$$\begin{aligned} \nabla_x^1 c(0)[s_4] + \nabla_x^2 c(0)[s_1, s_3] + \frac{1}{2} \nabla_x^2 c(0)[s_2]^2 + \frac{1}{2} \nabla_x^3 c(0)[s_1, s_1, s_2]^2 + \frac{1}{24} \nabla_x^4 c(0)[s_1]^4 \\ = 0 = e_2^T s_4 + e_1^T s_2 (e_1^T s_2 - 2\tau^2) + 2\tau e_1^T s_3 - 4\tau^2. \end{aligned}$$

The second, third and fourth of these conditions impose constraints on the values of $e_2^T s_2, e_2^T s_3$ and $e_2^T s_4$. In particular, the second implies that $e_2^T s_2 = -\tau^2$, which we may then substitute in (A.9) and deduce that

$$\begin{aligned} \nabla_x^4 \mu(0, t_\epsilon)[s_4] + \nabla_x^2 \mu(0, t_\epsilon)[s_1, s_3] + \frac{1}{2} \nabla_x^2 \mu(0, t_\epsilon)[s_2]^2 \\ + \frac{1}{2} \nabla_x^3 \mu(0, t_\epsilon)[s_1, s_1, s_2] + \frac{1}{24} \nabla_x^4 \mu(0, t_\epsilon)[s_1]^4 \\ = \frac{1}{2} \tau^4 + (\frac{1}{2} - \varepsilon) \tau^4 = (1 - \frac{1}{2} \varepsilon) \tau^4 \geq 0. \end{aligned} \tag{A.10}$$

We therefore obtain that, for all $\varepsilon \in (0, 1]$, x_* satisfies the necessary conditions of Theorem 2.1 with $q = 4$, except that $c(x_*) = \varepsilon$. However (A.5) shows that $\Lambda(x, y_\epsilon)$ is a polynomial of degree 4 with a global maximizer at the origin, independently of the value of ε . Letting ε tend to zero and using the fact that all quantities in the example depend continuously on this parameter then allows to conclude.

References

- [1] A. Anandkumar, R. Ge, Efficient approaches for escaping high-order saddle points in nonconvex optimization (2016). [arxiv:1602.05908](https://arxiv.org/abs/1602.05908).
- [2] E. Bergou, Y. Diouane, S. Gratton, On the use of the energy norm in trust-region and adaptive cubic regularization subproblems, April 2017.
- [3] W. Bian, X. Chen, Worst-case complexity of smoothing quadratic regularization methods for non-Lipschitzian optimization, SIAM J. Optim. 23 (3) (2013) 1718–1741.
- [4] W. Bian, X. Chen, Linearly constrained non-Lipschitzian optimization for image restoration, SIAM J. Imaging Sci. 8 (2015) 2294–2322.

- [5] W. Bian, X. Chen, Y. Ye, Complexity analysis of interior point algorithms for non-Lipschitz and nonconvex minimization, *Math. Program. A* 149 (2015) 301–327.
- [6] E.G. Birgin, J.L. Gardenghi, J.M. Martínez, S.A. Santos, Ph.L. Toint, Evaluation complexity for nonlinear constrained optimization using unscaled kkt conditions and high-order models, *SIAM J. Optim.* 26 (2) (2016) 951–967.
- [7] E.G. Birgin, J.L. Gardenghi, J.M. Martínez, S.A. Santos, Ph.L. Toint, Worst-case evaluation complexity for unconstrained nonlinear optimization using high-order regularized models, *Math. Program. A* 163 (1) (2017) 359–368.
- [8] J.F. Bonnans, R. Cominetti, A. Shapiro, Second order optimality conditions based on parabolic second order tangent sets, *SIAM J. Optim.* 9 (2) (1999) 466–492.
- [9] N. Boumal, P.-A. Absil, C. Cartis, Global rates of convergence for nonconvex optimization on manifolds (2016). [arxiv:1605.08101](https://arxiv.org/abs/1605.08101).
- [10] O.A. Brezhneva, A. Tret'yakov, Optimality conditions for degenerate extremum problems with equality constraints, *SIAM J. Control Optim.* 42 (2) (2003) 729–743.
- [11] O.A. Brezhneva, A. Tret'yakov, The p th-order optimality conditions for inequality constrained optimization problems, *Nonlinear Anal.* 63 (2005) 1357–1366.
- [12] C. Cartis, N.I.M. Gould, Ph.L. Toint, On the complexity of steepest descent, Newton's and regularized Newton's methods for nonconvex unconstrained optimization, *SIAM J. Optim.* 20 (6) (2010) 2833–2852.
- [13] C. Cartis, N.I.M. Gould, Ph.L. Toint, Adaptive cubic overestimation methods for unconstrained optimization. Part I: motivation, convergence and numerical results, *Math. Program. A* 127 (2) (2011) 245–295.
- [14] C. Cartis, N.I.M. Gould, Ph.L. Toint, Adaptive cubic overestimation methods for unconstrained optimization. Part II: worst-case function-evaluation complexity, *Math. Program. A* 130 (2) (2011) 295–319.
- [15] C. Cartis, N.I.M. Gould, Ph.L. Toint, An adaptive cubic regularization algorithm for nonconvex optimization with convex constraints and its function-evaluation complexity, *IMA J. Numer. Anal.* 32 (4) (2012) 1662–1695.
- [16] C. Cartis, N.I.M. Gould, Ph.L. Toint, On the complexity of finding first-order critical points in constrained nonlinear optimization, *Math. Program. A* 144 (1) (2013) 93–106.
- [17] C. Cartis, N.I.M. Gould, Ph.L. Toint, On the evaluation complexity of cubic regularization methods for potentially rank-deficient nonlinear least-squares problems and its relevance to constrained nonlinear optimization, *SIAM J. Optim.* 23 (3) (2013) 1553–1574.
- [18] C. Cartis, N.I.M. Gould, Ph.L. Toint, Evaluation complexity bounds for smooth constrained nonlinear optimization using scaled KKT conditions, high-order models and the χ criticality measure (2015). [arxiv:1705.04895](https://arxiv.org/abs/1705.04895).
- [19] C. Cartis, N.I.M. Gould, Ph.L. Toint, Evaluation complexity bounds for smooth constrained nonlinear optimization using scaled KKT conditions and high-order models, Technical Report naXys-(R1), Namur Center for Complex Systems (naXys), University of Namur, Namur, Belgium, 2015, pp. 11–2015.
- [20] C. Cartis, N.I.M. Gould, Ph.L. Toint, Improved worst-case evaluation complexity for potentially rank-deficient nonlinear least-Euclidean-norm problems using higher-order regularized models, Technical Report naXys-12-2015, Namur Center for Complex Systems (naXys), University of Namur, Namur, Belgium, 2015.
- [21] C. Cartis, N.I.M. Gould, Ph.L. Toint, On the evaluation complexity of constrained nonlinear least-squares and general constrained nonlinear optimization using second-order methods, *SIAM J. Numer. Anal.* 53 (2) (2015) 836–851.
- [22] C. Cartis, N.I.M. Gould, Ph.L. Toint, Second-order optimality and beyond: characterization and evaluation complexity in convexly constrained nonlinear optimization, *Found. Comput. Math.* 18 (5) (2018) 1073–1107.
- [23] C. Cartis, Ph.R. Sampaio, Ph.L. Toint, Worst-case complexity of first-order non-monotone gradient-related algorithms for unconstrained optimization, *Optimization* 64 (5) (2015) 1349–1361.
- [24] C. Cartis, K. Scheinberg, Global convergence rate analysis of unconstrained optimization methods based on probabilistic models, *Math. Program. A* (2017) <http://dx.doi.org/10.1007/s10107-017-1137-4> (in press).
- [25] X. Chen, Ph.L. Toint, H. Wang, Partially separable convexly-constrained optimization with non-Lipschitzian singularities and its complexity (2017). [arxiv:1704.06919](https://arxiv.org/abs/1704.06919).
- [26] R. Cominetti, Metric regularity, tangent sets and second-order optimality conditions, *Appl. Math. Optim.* 21 (1990) 265–287.
- [27] A.R. Conn, N.I.M. Gould, Ph.L. Toint, Trust-Region Methods, in: MPS-SIAM Series on Optimization, SIAM, Philadelphia, USA, 2000.
- [28] F.E. Curtis, D.P. Robinson, M. Samadi, Complexity analysis of a trust funnel algorithm for equality constrained optimization, Technical Report 16T-03, ISE/COR@L, LeHigh University, Bethlehem, PA, USA, 2017.
- [29] F.E. Curtis, D.P. Robinson, M. Samadi, A trust region algorithm with a worst-case iteration complexity of $O(\epsilon^{-3/2})$ for nonconvex optimization, *Math. Program. A* 162 (1) (2017) 1–32.
- [30] M. Dodangeh, L.N. Vicente, Z. Zhang, On the optimal order of worst case complexity of direct search, *Optim. Lett.* (2015) 1–10.
- [31] J.P. Dussault, Simple unified convergence proofs for the trust-region and a new ARC variant, Technical report, University of Sherbrooke, Sherbrooke, Canada, 2015.
- [32] R. Garmanjani, D. Júdice, L.N. Vicente, Trust-region methods without using derivatives: Worst case complexity and the non-smooth case, *SIAM J. Optim.* 26 (2016) 1987–2011.
- [33] D. Ge, X. Jiang, Y. Ye, A note on the complexity of l_p minimization, *Math. Program. A* 21 (2011) 1721–1739.
- [34] S. Ghadimi, G. Lan, Accelerated gradient methods for nonconvex nonlinear and stochastic programming, *Math. Program. A* 156 (1–2) (2016) 59–100.
- [35] G.N. Grapiglia, J. Yuan, Y. Yuan, On the convergence and worst-case complexity of trust-region and regularization methods for unconstrained optimization, *Math. Program. A* 152 (2015) 491–520.
- [36] G.N. Grapiglia, J. Yuan, Y. Yuan, Nonlinear stepsize control algorithms: Complexity bounds for first and second-order optimality, *J. Optim. Theory Appl.* 171 (2016) 971–997.

- [37] S. Gratton, C.W. Royer, L.N. Vicente, A decoupled first/second-order steps technique for nonconvex nonlinear unconstrained optimization with improved complexity bounds, Technical Report TR, Department of Mathematics, University of Coimbra, Coimbra, Portugal, 2017, pp. 17–21.
- [38] S. Gratton, C.W. Royer, L.N. Vicente, Z. Zhang, Direct search based on probabilistic descent, *SIAM J. Optim.* 25 (3) (2015) 1515–1541.
- [39] S. Gratton, A. Sartenaer, Ph.L. Toint, Recursive trust-region methods for multiscale nonlinear optimization, *SIAM J. Optim.* 19 (1) (2008) 414–444.
- [40] J.-B. Hiriart-Urruty, C. Lemaréchal, *Convex Analysis and Minimization Algorithms. Part 1: Fundamentals*, Springer Verlag, Heidelberg, Berlin, New York, 1993.
- [41] W. Hogan, Point-to-set maps in mathematical programming, *SIAM Rev.* 15 (3) (1973) 591–603.
- [42] F. Jarre, On Nesterov’s smooth Chebyshev-Rosenbrock function, *Optim. Methods Softw.* 28 (3) (2013) 478–500.
- [43] H. Kawasaki, An envelope-like effect of infinitely many inequality constraints on second order necessary conditions for minimization problems, *Math. Program.* 41 (1988) 73–96.
- [44] S. Łojasiewicz, *Ensembles semi-analytiques*, Technical report, Institut des Hautes Etudes Scientifiques, Bures-sur-Yvette, France, 1965.
- [45] S. Lu, Z. Wei, L. Li, A trust-region algorithm with adaptive cubic regularization methods for nonsmooth convex minimization, *Comput. Optim. Appl.* 51 (2012) 551–573.
- [46] J.M. Martínez, On high-order model regularization for constrained optimization, Technical report, Department of Applied Mathematics, IMECC-UNICAMP, Campinas, Brasil, 2017.
- [47] J.M. Martínez, M. Raydan, Cubic-regularization counterpart of a variable-norm trust-region method for unconstrained minimization, *J. Global Optim.* (2016) <http://dx.doi.org/10.1007/s10898-016-0475-8>.
- [48] Yu. Nesterov, G.N. Grapiglia, Globally convergent second-order schemes for minimizing twice-differentiable functions, Technical Report CORE Discussion paper 2016/28, CORE, Catholic University of Louvain, Louvain-la-Neuve, Belgium, 2016.
- [49] Yu. Nesterov, B.T. Polyak, Cubic regularization of Newton method and its global performance, *Math. Program. A* 108 (1) (2006) 177–205.
- [50] R.T. Rockafellar, *Convex Analysis*, Princeton University Press, Princeton, USA, 1970.
- [51] K. Scheinberg, X. Tang, Practical inexact proximal quasi-Newton method with global complexity analysis, *Math. Program. A* 160 (3) (2016).
- [52] K. Ueda, N. Yamashita, Convergence properties of the regularized Newton method for the unconstrained nonconvex optimization, *Appl. Math. Optim.* 62 (1) (2010) 27–46.
- [53] K. Ueda, N. Yamashita, On a global complexity bound of the levenberg-Marquardt method, *J. Optim. Theory Appl.* 147 (2010) 443–453.
- [54] L.N. Vicente, Worst case complexity of direct search, *EURO J. Comput. Optim.* 1 (2013) 143–153.