

RESEARCH OUTPUTS / RÉSULTATS DE RECHERCHE

Worst-case evaluation complexity of non-monotone gradient-related algorithms for unconstrained optimization

Cartis, C.; Rodrigues Sampaio, Phillipe; Toint, Ph L.

Published in:
Optimization

DOI:
[10.1080/02331934.2013.869809](https://doi.org/10.1080/02331934.2013.869809)

Publication date:
2015

Document Version
Early version, also known as pre-print

[Link to publication](#)

Citation for pulished version (HARVARD):
Cartis, C, Rodrigues Sampaio, P & Toint, PL 2015, 'Worst-case evaluation complexity of non-monotone gradient-related algorithms for unconstrained optimization', *Optimization*, vol. 64, no. 5, pp. 1349-1361. <https://doi.org/10.1080/02331934.2013.869809>

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal ?

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.



WORST-CASE EVALUATION COMPLEXITY OF
NON-MONOTONE GRADIENT-RELATED ALGORITHMS
FOR UNCONSTRAINED OPTIMIZATION

by C. Cartis¹, Ph. R. Sampaio² and Ph. L. Toint²

Report NAXYS-02-2013

12 March 2013



¹ University of Edinburgh, Edinburgh, EH9 3JZ, Scotland (UK)

² University of Namur, 61, rue de Bruxelles, B5000 Namur (Belgium)

<http://www.fundp.ac.be/sciences/naxys>

Worst-case evaluation complexity of non-monotone gradient-related algorithms for unconstrained optimization

C. Cartis*, Ph. R. Sampaio† and Ph. L. Toint‡

12 March 2013

Abstract

The worst-case evaluation complexity of finding an approximate first-order critical point using gradient-related non-monotone methods for smooth nonconvex and unconstrained problems is investigated. The analysis covers a practical linesearch implementation of these popular methods, allowing for an unknown number of evaluations of the objective function (and its gradient) per iteration. It is shown that this class of methods shares the known complexity properties of a simple steepest-descent scheme and that an approximate first-order critical point can be computed in at most $O(\epsilon^{-2})$ function and gradient evaluations, where $\epsilon > 0$ is the user-defined accuracy threshold on the gradient norm.

Keywords: Nonlinear optimization, evaluation complexity, worst-case analysis, linesearch algorithms, non-monotone methods.

1 Introduction

The worst-case evaluation complexity of optimization algorithms applied on nonlinear and potentially nonconvex problems has been studied in a sequence of recent papers, both for the unconstrained case (Nesterov, 2004, Gratton, Sartenaer and Toint, 2008, Nesterov and Polyak, 2006, Cartis, Gould and Toint, 2011a) and for the constrained one (Cartis, Gould and Toint, 2012a, 2013). Of particular interest here are the results of Nesterov (2004), page 29, in which this author analyzes the worst-case behaviour of the steepest-descent method for unconstrained minimization (both for exact and approximate linesearches) and shows that an approximate first-order critical point, that is a point at which the norm of the gradient of the objective function is less than $\epsilon > 0$, must be obtained in at most $O(\epsilon^{-2})$ iterations. Nesterov's analysis of the steepest-descent variants therefore effectively assumes that a single objective function value per iteration is computed, or at least that the number of such evaluations in the course of a single iteration is bounded. His bounds thus specify iteration-complexity rather than evaluation complexity. At variance, more

*School of Mathematics, University of Edinburgh, The King's Buildings, Edinburgh, EH9 3JZ, Scotland, UK. Email: coralia.cartis@ed.ac.uk

†Namur Center for Complex Systems (naXys) and Department of Mathematics, FUNDP-University of Namur, 61, rue de Bruxelles, B-5000 Namur, Belgium. Email: philippe.toint@unamur.be

‡Namur Center for Complex Systems (naXys) and Department of Mathematics, University of Namur, 61, rue de Bruxelles, B-5000 Namur, Belgium. Email: philippe.toint@unamur.be

typical implementations use a linesearch (which makes no explicit use of the Lipschitz constant) to compute a suitable steplength, with the possible drawback that an unknown number of additional function evaluations may be required during the course of a single iteration. The question of the worst-case objective-function evaluation complexity of linesearch implementations of this type has not yet been considered specifically. Interestingly, a worst-case complexity analysis is available for other first-order algorithms, such as first-order trust-region methods (Gratton et al., 2008) and first-order regularization algorithms (Cartis, Gould and Toint, 2011b).

In parallel, it has long been known that “gradient related” minimization methods share a number of their convergence properties with the steepest-descent algorithm (see Ortega and Rheinboldt, 1970 for an early reference). In these methods, a linesearch is performed along a direction whose angle with the negative gradient is bounded away from orthogonality. This class covers a wide range of practical algorithms, including for instance variable-metric techniques or finite-difference schemes when Hessian approximations have bounded conditioning (see Nocedal and Wright, 1999, page 40, for instance). Despite their close connection with steepest descent, their worst-case analysis remains so far an open question.

Standard linesearch methods are usually defined in a way which ensures monotonically decreasing objective-function values as the iterations proceed. However, “non-monotone” generalizations of these algorithms, where this monotonicity property is abandoned, have gained respect in practice because of their often better performance. We refer the reader to Grippo, Lampariello and Lucidi 1986, 1989, or Toint (1996) for more details on these methods. Again, the worst-case performance of this interesting class of algorithms is so far unexplored.

The purpose of this paper is to bring together these three questions (standard linesearch, gradient-related directions and non-monotonicity) and to provide an analysis which covers them all. We therefore consider non-monotone gradient-related linesearch optimization methods and show that, as for steepest-descent, their objective-function evaluation complexity is $O(\epsilon^{-2})$. Note that standard monotone variants are also covered by this analysis.

Section 2 states the problem and describes the class of algorithms considered, while Section 3 provides an upper bound on their worst-case evaluation complexity. Some comments and perspectives are finally presented in Section 4.

2 The problem and algorithm

We consider the nonlinear and possibly nonconvex smooth unconstrained minimization problem

$$\min_{x \in \mathbb{R}^n} f(x) \tag{2.1}$$

for which we assume the following.

AF0 $f(x)$ is bounded below on \mathbb{R}^n , that is there exists a constant⁽¹⁾ $\kappa_{\text{lb}f}$ such that, for all $x \in \mathbb{R}^n$, $f(x) \geq \kappa_{\text{lb}f}$.

⁽¹⁾“lb” stands for “lower bound on the objective function”.

AF1 $f(x)$ is continuously differentiable on \mathbb{R}^n .

As stated in the introduction, we consider a class of algorithm in which the search directions are “gradient-related” (see Ortega and Rheinboldt, 1970, and Bertsekas, 2008, page 35). This terminology means that, at iteration k , an approximate unidimensional minimization of the objective function is performed along a direction d_k whose angle with the steepest descent is controlled by the condition

$$\langle g_k, d_k \rangle \leq -\kappa_1 \|g_k\|^2 \text{ and } \|d_k\| \leq \kappa_2 \|g_k\|, \quad (2.2)$$

where $g_k \stackrel{\text{def}}{=} g(x_k) \stackrel{\text{def}}{=} \nabla_x f(x_k)$, x_k is the k -th iterate, $\langle \cdot, \cdot \rangle$ and $\|\cdot\|$ are the Euclidean inner product and norm, respectively, and κ_1 and κ_2 are positive constants independent of k .

Once the direction is fixed, it is then used in a non-monotone linesearch. We choose here a Goldstein-Armijo variant (see Grippo et al., 1986, or Nocedal and Wright, 1999, pages 33-37), in which a stepsize t_k (yielding a new iterate $x_{k+1} = x_k + t_k d_k$) is accepted whenever the conditions

$$f(x_k + t_k d_k) \leq \max_{0 \leq j \leq M} [f(x_{k-j})] + \alpha t_k \langle g_k, d_k \rangle, \quad (2.3)$$

and

$$f(x_k + t_k d_k) \geq \max_{0 \leq j \leq M} [f(x_{k-j})] + \beta t_k \langle g_k, d_k \rangle \quad (2.4)$$

hold, where $M \geq 0$, $\alpha \in (0, 1)$ and $\beta \in (\alpha, 1)$ are constants independent of k , and where, by convention, $x_{-M} = \dots x_{-1} = x_0$. Note that $M = 0$ corresponds to the monotone case.

The class of algorithms of interest may now be stated formally as Algorithm 2.1 on the following page.

Note that the successive phases of the Goldstein-Armijo technique are apparent in the algorithm’s description: a bracket containing the desired step is first identified by bactracking (Step 4) or look-ahead (Step 5), and the final step is then computed by bisection (Step 6).

3 Worst-case evaluation complexity analysis

We now analyze the worst-case behaviour of Algorithm 2.1. A first step in this analysis is to specify our assumptions.

AF2 $g(x)$ is Lipschitz continuous on \mathbb{R}^n , that is there exists a constant $L_g > 0$ such that, for all $x, y \in \mathbb{R}^n$,

$$\|g(x) - g(y)\| \leq L_g \|x - y\|.$$

The first simple but crucial property that can be deduced from these assumptions is that the stepsize is bounded below by a constant inversely proportional to the Lipschitz constant L_g .

Algorithm 2.1: A gradient-related non-monotone linesearch algorithm.

Step 0: Initialization. An initial point x_0 is given, as well as an accuracy level $\epsilon > 0$. The constants t_{ini} , M , α and β are also given, satisfying $t_{\text{ini}} > 0$, $M \geq 0$ and $0 < \alpha < \beta < 1$. Compute $f(x_0)$, g_0 and set $k = 0$.

Step 1: Test for termination. If $\|g_k\| \leq \epsilon$, terminate.

Step 2: Select a search direction. Choose d_k such that (2.2) holds.

Step 3: Linesearch: test initial stepsize.

1. Set $t_k = t_{\text{ini}} > 0$, $t_{\text{low}} = 0$ and compute $f(x_k + t_k d_k)$.
2. If (2.3) fails, go to Step 4.
3. If (2.4) fails, go to Step 5.
4. Else go to Step 7.

Step 4: Linesearch: backtracking.

1. While (2.3) fails, set $t_{\text{up}} \leftarrow t_k$, $t_k \leftarrow \frac{1}{2}t_k$ and compute $f(x_k + t_k d_k)$.
2. If (2.4) holds, go to Step 7, or set $t_{\text{low}} \leftarrow t_k$ and go to Step 6 otherwise.

Step 5: Linesearch: look ahead.

1. While (2.4) fails, set $t_{\text{low}} \leftarrow t_k$, $t_k \leftarrow 2t_k$ and compute $f(x_k + t_k d_k)$.
2. If (2.3) holds, go to Step 7, or set $t_{\text{up}} \leftarrow t_k$ and go to Step 6 otherwise.

Step 6: Linesearch: bisection inside bracket.

1. Set $t_k \leftarrow \frac{1}{2}(t_{\text{low}} + t_{\text{up}})$ and compute $f(x_k + t_k d_k)$.
2. If (2.3) fails, set $t_{\text{up}} \leftarrow t_k$ and return to Step 6.
3. If (2.4) fails, set $t_{\text{low}} \leftarrow t_k$ and return to Step 6.

Step 7: Compute the new iterate and gradient. Set $x_{k+1} = x_k + t_k d_k$ and compute $g_{k+1} = g(x_{k+1})$. Increment k by one and return to Step 1.

Lemma 3.1 Suppose that AF0–AF2 hold. Then any value of $t > 0$ such that (2.4) holds for $t_k = t$ also satisfies the inequality

$$t \geq \frac{2(1-\beta)\kappa_1}{L_g \kappa_2^2}. \quad (3.1)$$

Proof. We successively use the mean value theorem, the Cauchy-Schwarz inequality and AF2

to obtain that

$$\begin{aligned}
f(x_k + td_k) &= f(x_k) + t\langle g_k, d_k \rangle + \int_0^1 \langle g(x_k + \tau td_k) - g_k, td_k \rangle d\tau \\
&\leq f(x_k) + t\langle g_k, d_k \rangle + t\|d_k\| \int_0^1 \|g(x_k + \tau td_k) - g_k\| d\tau \\
&\leq f(x_k) + t\langle g_k, d_k \rangle + \frac{1}{2}t^2 L_g \|d_k\|^2 \\
&\leq \max_{0 \leq j \leq M} [f(x_{k-j})] + t\langle g_k, d_k \rangle + \frac{1}{2}t^2 L_g \|d_k\|^2.
\end{aligned} \tag{3.2}$$

Combining this relation with (2.4) and (2.2), we have that

$$t \geq \frac{2\langle g_k, d_k \rangle (\beta - 1)}{L_g \|d_k\|^2} \geq \frac{2(1 - \beta)\kappa_1 \|g_k\|^2}{L_g \|g_k\|^2 \kappa_2^2} = \frac{2(1 - \beta)\kappa_1}{L_g \kappa_2^2}. \tag{3.3}$$

□

We now prove that there is a finite and non-empty interval of acceptable stepsizes.

Lemma 3.2 Suppose that AF0-AF1 hold and that $g_k \neq 0$. Then there exists an interval $[t_k^\beta, t_k^\alpha]$ such that

$$0 < t_k^\beta < t_k^\alpha < +\infty \tag{3.4}$$

and (2.3)-(2.4) hold for every value of $t_k \in [t_k^\beta, t_k^\alpha]$.

Proof. Observe first that the slope of $f(x_k + td_k)$ is steeper than that of the straight lines $f(x_k) + \alpha t\langle g_k, d_k \rangle$ and $f(x_k) + \beta t\langle g_k, d_k \rangle$, ($t \geq 0$), since $\alpha < 1$ and $\beta < 1$. Thus, for all $t > 0$ sufficiently small,

$$f(x_k + td_k) < f(x_k) + \alpha t\langle g_k, d_k \rangle \leq \max_{0 \leq j \leq M} f(x_{k-j}) + \alpha t\langle g_k, d_k \rangle \tag{3.5}$$

and

$$f(x_k + td_k) < f(x_k) + \beta t\langle g_k, d_k \rangle \leq \max_{0 \leq j \leq M} f(x_{k-j}) + \beta t\langle g_k, d_k \rangle. \tag{3.6}$$

It follows from (3.5) that (2.3) holds for all t_k sufficiently small. Furthermore, (2.3) does not hold in the limit as $t_k = t \rightarrow \infty$ since $f(x_k) + \alpha t\langle g_k, d_k \rangle \leq f(x_k) - \alpha t\kappa_1 \|g_k\|^2 \rightarrow -\infty$ (because of (2.2)), while $f(x_k + td_k) \geq \kappa_{\text{lb}} f$ for all t due to AF0. Thus there exists a value $0 < t_k^\alpha < \infty$ such that

$$f(x_k + t_k^\alpha d_k) = \max_{0 \leq j \leq M} f(x_{k-j}) + \alpha t_k^\alpha \langle g_k, d_k \rangle. \tag{3.7}$$

For simplicity, let us choose the smallest t_k^α that satisfies (3.7) so that (3.5) holds for all $t \in$

$(0, t_k^\alpha)$. Since $\alpha < \beta < 1$, we note that

$$\max_{0 \leq j \leq M} f(x_{k-j}) + \beta t \langle g_k, d_k \rangle < \max_{0 \leq j \leq M} f(x_{k-j}) + \alpha t \langle g_k, d_k \rangle \quad \text{for all } t > 0.$$

Letting $t = t_k^\alpha$ in this inequality and using (3.7), we deduce that (2.4) must continue to hold for $0 < t_k = t < t_k^\alpha$ sufficiently close to t_k^α . However, (3.6) implies that (2.4) must fail for sufficiently small $t > 0$, and using again AF1, we conclude that there exists $0 < t_k^\beta < t_k^\alpha$ such that

$$f(x_k + t_k^\beta d_k) = \max_{0 \leq j \leq M} f(x_{k-j}) + \beta t_k^\beta \langle g_k, d_k \rangle, \quad (3.8)$$

and (2.4) holds for all $t_k \in [t_k^\beta, t_k^\alpha]$. (Clearly, t_k^β must be distinct from $t_k^\alpha < \infty$ due to (3.7), (3.8) and $\alpha < \beta$.) This concludes the proof since (2.3) holds for t_k in the same interval due to (3.5) and the definition of t_k^α . \square

Having proved the existence of an interval of acceptable stepsizes, we now verify that the measure of this interval is bounded below by some positive constant.

Lemma 3.3 Suppose that AF0-AF2 hold, and define t_k^α and t_k^β to be any solutions of (3.7) and (3.8), respectively, such that (2.3) and (2.4) hold for each $t \in [t_k^\beta, t_k^\alpha]$. Then the interval $[t_k^\beta, t_k^\alpha]$ has a strictly positive measure in the sense that there exists a constant $\kappa_{\text{int}} > 0$ only depending on $\alpha, \beta, \kappa_1, \kappa_2$ and L_g such that

$$t_k^\alpha - t_k^\beta \geq \kappa_{\text{int}}. \quad (3.9)$$

Proof. Assume first that $f(x_k + t_k^\alpha d_k) \leq f(x_k + t_k^\beta d_k)$. Then (3.7) and (3.8) imply that $\alpha t_k^\alpha > \beta t_k^\beta$, and so, using also Lemma 3.1,

$$t_k^\alpha - t_k^\beta \geq \frac{\beta - \alpha}{\alpha} t_k^\beta \geq \frac{2(\beta - \alpha)(1 - \beta)\kappa_1}{\alpha L_g \kappa_2^2}. \quad (3.10)$$

Suppose now that $f(x_k + t_k^\alpha d_k) > f(x_k + t_k^\beta d_k)$. Applying the mean value theorem to $f(x + td)$ on $[t_k^\beta, t_k^\alpha]$ yields that

$$\begin{aligned} f(x_k + t_k^\alpha d_k) - f(x_k + t_k^\beta d_k) &= (t_k^\alpha - t_k^\beta) \langle g(x_k + t_\xi d_k), d_k \rangle \\ &\leq (t_k^\alpha - t_k^\beta) \|g(x_k + t_\xi d_k)\| \|d_k\| \\ &\leq (t_k^\alpha - t_k^\beta) \kappa_2 \|g(x_k + t_\xi d_k)\| \|g_k\|, \end{aligned}$$

where $t_\xi \in (t_k^\beta, t_k^\alpha)$ and where the first inequality follows from the Cauchy-Schwarz inequality and the second from (2.2). Furthermore, the Lipschitz continuity of g (AF2), (2.2) and the

bound $t_\xi < t_k^\alpha$ give that

$$\begin{aligned} \|g(x_k + t_\xi d_k)\| &\leq \|g(x_k + t_\xi d_k) - g_k\| + \|g_k\| \\ &\leq L_g t_\xi \|d_k\| + \|g_k\| \\ &\leq L_g t_\xi \kappa_2 \|g_k\| + \|g_k\| \\ &\leq (L_g t_k^\alpha \kappa_2 + 1) \|g_k\|. \end{aligned}$$

Thus

$$f(x_k + t_k^\alpha d_k) - f(x_k + t_k^\beta d_k) \leq \kappa_2 (t_k^\alpha - t_k^\beta) (L_g t_k^\alpha \kappa_2 + 1) \|g_k\|^2. \quad (3.11)$$

The definition of t_k^α and t_k^β in Lemma 3.2 then give that (3.7) and (3.8) both hold, and so

$$\begin{aligned} f(x_k + t_k^\alpha d_k) - f(x_k + t_k^\beta d_k) &= \alpha t_k^\alpha \langle g_k, d_k \rangle - \beta t_k^\beta \langle g_k, d_k \rangle \\ &= (\alpha t_k^\alpha - \beta t_k^\beta) \langle g_k, d_k \rangle \\ &\geq (\beta t_k^\beta - \alpha t_k^\alpha) \kappa_1 \|g_k\|^2, \end{aligned} \quad (3.12)$$

where we again used the Cauchy-Schwartz inequality and (2.2) to deduce the last inequality. From (3.11), we now deduce that

$$(\beta t_k^\beta - \alpha t_k^\alpha) \kappa_1 \leq \kappa_2 (t_k^\alpha - t_k^\beta) (L_g t_k^\alpha \kappa_2 + 1). \quad (3.13)$$

This inequality is equivalent to

$$\kappa_2^2 L_g (t_k^\alpha)^2 + (\kappa_2 + \alpha \kappa_1 - \kappa_2^2 L_g t_k^\beta) t_k^\alpha - (\kappa_2 + \beta \kappa_1) t_k^\beta \geq 0, \quad (3.14)$$

and so, since $t_k^\alpha > 0$, we deduce that

$$t_k^\alpha \geq \frac{\kappa_2^2 L_g t_k^\beta - \kappa_2 - \alpha \kappa_1 + \sqrt{(\kappa_2 + \alpha \kappa_1 - \kappa_2^2 L_g t_k^\beta)^2 + 4(\kappa_2 + \beta \kappa_1) \kappa_2^2 L_g t_k^\beta}}{2\kappa_2^2 L_g}, \quad (3.15)$$

and therefore that

$$\begin{aligned} 2\kappa_2^2 L_g (t_k^\alpha - t_k^\beta) &\geq -(\kappa_2 + \alpha \kappa_1 + \kappa_2^2 L_g t_k^\beta) + \sqrt{(\kappa_2 + \alpha \kappa_1 - \kappa_2^2 L_g t_k^\beta)^2 + 4(\kappa_2 + \beta \kappa_1) \kappa_2^2 L_g t_k^\beta} \\ &= \frac{-(\kappa_2 + \alpha \kappa_1 + \kappa_2^2 L_g t_k^\beta)^2 + (\kappa_2 + \alpha \kappa_1 - \kappa_2^2 L_g t_k^\beta)^2 + 4(\kappa_2 + \beta \kappa_1) \kappa_2^2 L_g t_k^\beta}{\kappa_2 + \alpha \kappa_1 + \kappa_2^2 L_g t_k^\beta + \sqrt{(\kappa_2 + \alpha \kappa_1 - \kappa_2^2 L_g t_k^\beta)^2 + 4(\kappa_2 + \beta \kappa_1) \kappa_2^2 L_g t_k^\beta}} \\ &= \frac{4(\beta - \alpha) \kappa_1 \kappa_2^2 L_g t_k^\beta}{\kappa_2 + \alpha \kappa_1 + \kappa_2^2 L_g t_k^\beta + \sqrt{(\kappa_2 + \alpha \kappa_1 - \kappa_2^2 L_g t_k^\beta)^2 + 4(\kappa_2 + \beta \kappa_1) \kappa_2^2 L_g t_k^\beta}}. \end{aligned} \quad (3.16)$$

As a consequence, we obtain that

$$\frac{(t_k^\alpha - t_k^\beta)}{2(\beta - \alpha) \kappa_1} \geq \frac{t_k^\beta}{\kappa_2 + \alpha \kappa_1 + \kappa_2^2 L_g t_k^\beta + \sqrt{(\kappa_2 + \alpha \kappa_1 - \kappa_2^2 L_g t_k^\beta)^2 + 4(\kappa_2 + \beta \kappa_1) \kappa_2^2 L_g t_k^\beta}} \stackrel{\text{def}}{=} E(t_k^\beta). \quad (3.17)$$

Defining $S(t_k^\beta) \stackrel{\text{def}}{=} \sqrt{(\kappa_2 + \alpha\kappa_1 - \kappa_2^2 L_g t_k^\beta)^2 + 4(\kappa_2 + \beta\kappa_1)\kappa_2^2 L_g t_k^\beta}$, differentiating $E(t_k^\beta)$ with respect to t_k^β then gives that

$$\begin{aligned}
E'(t_k^\beta) &= \\
& \frac{\kappa_2 + \alpha\kappa_1 + \kappa_2^2 L_g t_k^\beta + S(t_k^\beta) - t_k^\beta \left[\kappa_2^2 L_g + \frac{-(\kappa_2 + \alpha\kappa_1 - \kappa_2^2 L_g t_k^\beta)\kappa_2^2 L_g + 2(\kappa_2 + \beta\kappa_1)\kappa_2^2 L_g}{S(t_k^\beta)} \right]}{\left[\kappa_2 + \alpha\kappa_1 + \kappa_2^2 L_g t_k^\beta + S(t_k^\beta) \right]^2} \\
&= \frac{(\kappa_2 + \alpha\kappa_1)S(t_k^\beta) + (\kappa_2 + \alpha\kappa_1 - \kappa_2^2 L_g t_k^\beta)^2 + (\kappa_2 + \alpha\kappa_1 - \kappa_2^2 L_g t_k^\beta)\kappa_2^2 L_g t_k^\beta + 2(\kappa_2 + \beta\kappa_1)\kappa_2^2 L_g t_k^\beta}{\left[\kappa_2 + \alpha\kappa_1 + \kappa_2^2 L_g t_k^\beta + S(t_k^\beta) \right]^2 S(t_k^\beta)} \\
&= \frac{(\kappa_2 + \alpha\kappa_1)[S(t_k^\beta) + \kappa_2 + \alpha\kappa_1 - \kappa_2^2 L_g t_k^\beta] + 2(\kappa_2 + \beta\kappa_1)\kappa_2^2 L_g t_k^\beta}{\left[\kappa_2 + \alpha\kappa_1 + \kappa_2^2 L_g t_k^\beta + S(t_k^\beta) \right]^2 S(t_k^\beta)}.
\end{aligned} \tag{3.18}$$

It then follows that

$$E'(t_k^\beta) > 0 \text{ for all } t_k^\beta > 0$$

since

$$S(t_k^\beta) + \kappa_2 + \alpha\kappa_1 - \kappa_2^2 L_g t_k^\beta > |\kappa_2 + \alpha\kappa_1 - \kappa_2^2 L_g t_k^\beta| + \kappa_2 + \alpha\kappa_1 - \kappa_2^2 L_g t_k^\beta \geq 0$$

and each constant and variable in $E'(t_k^\beta)$ is positive. Thus $E(t_k^\beta)$ is increasing as a function of t_k^β , and we obtain, because of Lemma 3.1 and the fact that (2.4) holds at t_k^β by construction, that

$$E(t_k^\beta) \geq E\left(\frac{2(1-\beta)\kappa_1}{L_g \kappa_2^2}\right),$$

and we finally deduce from (3.18) that

$$t_k^\alpha - t_k^\beta \geq 2(\beta - \alpha)\kappa_1 E\left(\frac{2(1-\beta)\kappa_1}{L_g \kappa_2^2}\right).$$

Combining this with (3.10), we deduce that (3.9) holds with

$$\kappa_{\text{int}} \stackrel{\text{def}}{=} 2(\beta - \alpha)\kappa_1 \min\left[\frac{(1-\beta)}{\alpha L_g \kappa_2^2}, E\left(\frac{2(1-\beta)\kappa_1}{L_g \kappa_2^2}\right)\right],$$

where this lower bound only depends on α , β , κ_1 , κ_2 and L_g , as desired. \square

We now turn to estimating the worst-case evaluation complexity of Algorithm 2.1 for the task of finding an ϵ -first-order critical point.

Theorem 3.4 Suppose that AF0-AF2 hold. Then there exists a constant $\kappa_{\text{GNL}} > 0$ such that, for any $\epsilon \in (0, 1)$, Algorithm 2.1 needs at most

$$\left\lceil \frac{\kappa_{\text{GNL}}(f(x_0) - \kappa_{\text{lbF}})}{\epsilon^2} + M \right\rceil$$

to produce an iterate x_k such that $\|g_k\| \leq \epsilon$, where κ_{lbF} and M are defined in AF0 and (2.3)-(2.4), respectively, and where

$$\kappa_{\text{GNL}} \stackrel{\text{def}}{=} \frac{(M+1)}{\alpha\kappa_1} \max \left[\frac{L_g\kappa_2^2 \max[n_1, n_2]}{2(1-\beta)\kappa_1}, \frac{2(n_2+1)}{t_{\text{ini}}} \right]$$

with

$$n_1 \stackrel{\text{def}}{=} \left\lceil \log_2 \left(\frac{(1-\beta)\kappa_1}{t_{\text{ini}}L_g\kappa_2^2} \right) \right\rceil \quad \text{and} \quad n_2 \stackrel{\text{def}}{=} \left\lceil \log_2 \left(\frac{\kappa_{\text{int}}}{t_{\text{ini}}} \right) \right\rceil.$$

Proof. The proof proceeds by first establishing the minimum achieved decrease in the objective function between iterate x_{k+1} and its “predecessor” $x_{\pi(k+1)}$, where

$$\pi(k+1) = k - \arg \max_{0 \leq j \leq M} f(x_{k-j}) \quad (3.19)$$

when using Algorithm 2.1.

- Assume first that both (2.3) and (2.4) hold for $t_k = t_{\text{ini}}$ (in Step 3). Then we obtain a decrease

$$f(x_{\pi(k+1)}) - f(x_{k+1}) \geq -\alpha t_{\text{ini}} \langle g_k, d_k \rangle \geq \alpha t_{\text{ini}} \kappa_1 \|g_k\|^2, \quad (3.20)$$

because of (2.3) and (2.2), and this decrease is obtained for a single additional function evaluation.

- Assume now that (2.3) fails at Step 3.2, and Step 4 is therefore entered. Assume furthermore that $j_3 \geq 1$ backtracking steps are performed in Step 4.1. The j_3 is the smallest non-negative integer such that (2.3) holds for $t_k = t_{\text{ini}}2^{-j_3}$, which means that j_3 is the largest integer for which this inequality is violated for $t = t_{\text{ini}}2^{-j_3+1}$. Because $\alpha < \beta$, we deduce that (2.4) must hold for this value of t_k . Using now Lemma 3.1, we obtain that

$$t = 2^{-j_3+1}t_{\text{ini}} \geq \frac{2(1-\beta)\kappa_1}{L_g\kappa_2^2},$$

which in turn implies that

$$j_3 \leq \left\lceil \log_2 \left(\frac{(1-\beta)\kappa_1}{t_{\text{ini}}L_g\kappa_2^2} \right) \right\rceil \stackrel{\text{def}}{=} n_1. \quad (3.21)$$

Step 4 therefore requires at most n_1 function evaluations. If the linesearch is terminated in Step 4.2 (i.e., branching occurs to Step 7), we obtain a decrease

$$f(x_{\pi(k+1)}) - f(x_{k+1}) \geq -\alpha t_k \langle g_k, d_k \rangle \geq \alpha \frac{2(1-\beta)\kappa_1}{L_g \kappa_2^2} \kappa_1 \|g_k\|^2, \quad (3.22)$$

where we used (2.3), (2.4), the Cauchy-Schwarz inequality, (2.2) and Lemma 3.1 successively.

- If the linesearch is not terminated in Step 4, Step 6 must be entered, with a bracket $[t_{\text{low}}, t_{\text{up}}]$ where

$$t_{\text{low}} = 2^{-j_3} t_{\text{ini}} = \frac{1}{2} t_{\text{up}}.$$

Thus

$$t_{\text{up}} - t_{\text{low}} = \frac{1}{2} 2^{-j_3+1} t_{\text{ini}} = 2^{-j_3} t_{\text{ini}}.$$

We know from Lemma 3.3 that the length of the admissible interval is at least equal to $\kappa_{\text{int}} > 0$, where this constant only depends on α , β and L_g . Thus the number $j_4 \geq 1$ of bisection (and function evaluations) within Step 6 is bounded above by the smallest integer such that

$$2^{-j_4} (t_{\text{up}} - t_{\text{low}}) = 2^{-j_4} 2^{-j_3} t_{\text{ini}} \geq \kappa_{\text{int}},$$

which then yields that the total number of function evaluations in Step 4 and 6 is bounded by

$$j_3 + j_4 \leq \left\lceil \log_2 \left(\frac{\kappa_{\text{int}}}{t_{\text{ini}}} \right) \right\rceil \stackrel{\text{def}}{=} n_2.$$

If we know compute the decrease obtained, we deduce, again from (2.3), (2.4) and Lemma 3.1, that (3.22) also holds in this case.

- Assume now that (2.4) fails in Step 3.3, and thus that Step 5 is entered. Assume furthermore that $j_2 \geq 1$ doubling of t_k (and j_2 function evaluations) occur in Step 5.1 (we know that j_2 is finite because of (3.4)). If the linesearch is terminated in Step 5.2 (i.e., branching to Step 7 occurs), we obtain that the function decrease obtained is bounded below by

$$f(x_{\pi(k+1)}) - f(x_{k+1}) \geq -\alpha t_k \langle g_k, d_k \rangle \geq \alpha 2^{j_2} t_{\text{ini}} \kappa_1 \|g_k\|^2.$$

- The final case is when Step 6 is entered after Step 5, in which case the initial bracket for Step 6 is given by $[t_{\text{low}}, t_{\text{up}}]$ where

$$t_{\text{low}} = 2^{j_2-1} t_{\text{ini}} = \frac{1}{2} t_{\text{up}}.$$

Thus

$$t_{\text{up}} - t_{\text{low}} = \frac{1}{2} 2^{j_2} t_{\text{ini}} = 2^{j_2-1} t_{\text{ini}}.$$

Just as in the case where Step 6 is entered after Step 4, we now deduce that the number

j_4 of bisections and function evaluations needed to reduce this bracket to the minimum possible value κ_{int} is limited by the inequality

$$2^{-j_4} 2^{j_2-1} t_{\text{ini}} = 2^{-j_4} (t_{\text{up}} - t_{\text{low}}) \geq \kappa_{\text{int}},$$

yielding a maximum number of bisection (and function evaluation) in Step 6 bounded by

$$j_4 \leq j_2 - 1 + \left\lceil \log_2 \left(\frac{\kappa_{\text{int}}}{t_{\text{ini}}} \right) \right\rceil = j_2 - 1 + n_2 \leq j_2(n_2 + 1).$$

In this final case, since $t_k \geq t_{\text{low}} = 2^{j_2-1} t_{\text{ini}}$ and (2.3) holds at t_k , the function decrease is bounded below by

$$f(x_{\pi(k+1)}) - f(x_{k+1}) \geq -\alpha t_k \langle g_k, d_k \rangle \geq \alpha 2^{j_2-1} t_{\text{ini}} \kappa_1 \|g_k\|^2.$$

Gathering all cases together, we see that function decrease per function evaluation is given, in the worst case, by

$$\min \left[\frac{t_{\text{ini}}}{1}, \frac{2(1-\beta)\kappa_1}{L_g \kappa_2^2 n_1}, \frac{2(1-\beta)\kappa_1}{L_g \kappa_2^2 n_2}, \frac{2^{j_2} t_{\text{ini}}}{j_2}, \frac{2^{j_2-1} t_{\text{ini}}}{2j_2(n_2+1)} \right] \alpha \kappa_1 \|g_k\|^2, \quad (3.23)$$

where, by construction, n_1 and n_2 only depend on α , β , κ_1 , κ_2 , L_g and t_{ini} . Noting that, for $j_2 \geq 1$,

$$\frac{2^{j_2}}{j_2} \geq 2 \quad \text{and} \quad \frac{2^{j_2-1}}{2j_2} \geq \frac{1}{2}$$

and defining

$$\kappa_{\text{decr}} \stackrel{\text{def}}{=} \alpha \kappa_1 \min \left[\frac{2(1-\beta)\kappa_1}{L_g \kappa_2^2 \max[n_1, n_2]}, \frac{t_{\text{ini}}}{2(n_2+1)} \right],$$

we therefore deduce from (3.23) that, as long as the algorithm does not terminate (i.e., as long as $\|g_k\| \geq \epsilon$)

$$f(x_{\pi(k+1)}) - f(x_{k+1}) \geq \kappa_{\text{decr}} \|g_k\|^2 \geq \kappa_{\text{decr}} \epsilon^2.$$

Tracing back the predecessors of iterate x_{k+1} up to x_0 and denoting the composition of j instances of the predecessor operator $\pi(\cdot)$ by $\pi^j(\cdot)$, we also deduce that

$$f(x_{\pi^{j+1}(k+1)}) - f(x_{\pi^j(k+1)}) \geq \kappa_{\text{decr}} \epsilon^2 \quad (3.24)$$

for all $j = 0, \dots, p_k$ for p_k such that $x_{\pi^{p_k}(k+1)} = x_0$ and where, by convention, $\pi^0(k+1) \stackrel{\text{def}}{=} k+1$. Now the definition of $\pi(\cdot)$ in (3.19) implies that, for all ℓ ,

$$0 \leq \ell + 1 - \pi(\ell + 1) \leq M + 1, \quad (3.25)$$

and we have, using AF0, that

$$f(x_0) - \kappa_{\text{lb}} \geq f(x_0) - f(x_{k+1}) = \sum_{j=0}^{P_k} [f(x_{\pi^{j+1}(k+1)}) - f(x_{\pi^j(k+1)})]. \quad (3.26)$$

Using (3.25), we obtain that the sum in the right-side of this expression contains at least

$$\left\lfloor \frac{k+1}{M+1} \right\rfloor$$

terms. Substituting then (3.24) for each term, (3.26) gives that

$$f(x_0) - \kappa_{\text{lb}} \geq \left\lfloor \frac{k+1}{M+1} \right\rfloor \kappa_{\text{decr}} \epsilon^2 \geq \left(\frac{k+1}{M+1} - 1 \right) \kappa_{\text{decr}} \epsilon^2 = \frac{k-M}{M+1} \kappa_{\text{decr}} \epsilon^2.$$

As a consequence, we obtain that the total number of function evaluations in Algorithm 2.1 is bounded above by

$$\left\lceil \frac{(M+1)(f(x_0) - \kappa_{\text{lb}})}{\kappa_{\text{decr}} \epsilon^2} + M \right\rceil$$

yielding the desired conclusion with $\kappa_{\text{GNL}} = (M+1)/\kappa_{\text{decr}}$. \square

4 Conclusions and perspectives

We have shown that gradient-related methods using a non-monotone (and monotone) linesearch will find an ϵ -approximate first-order critical point of a smooth function with Lipschitz gradient in $O(\epsilon^{-2})$ function and gradient evaluations at most. Their worst-case behaviour is therefore, up to a factor, equivalent to that of a simple monotone pure steepest-descent algorithm, albeit their practical performance is often superior (see Toint, 1996). Moreover, it results from Cartis, Gould and Toint (2010), that this bound is sharp.

In the same line of investigation, Cartis, Gould and Toint (2012b) show that the same complexity order is obtained for the steepest-descent method with exact linesearch and that it is sharp. One may expect that this result can be extended to the gradient-related algorithms analyzed in the present note, although the construction of an example illustrating the sharpness of the complexity bound is likely to be challenging without monotonicity.

References

- D. P. Bertsekas. *Nonlinear Programming*. Athena Scientific, Belmont, Massachusetts, USA, 2008.
- C. Cartis, N. I. M. Gould, and Ph. L. Toint. On the complexity of steepest descent, Newton's and regularized Newton's methods for nonconvex unconstrained optimization. *SIAM Journal on Optimization*, **20**(6), 2833–2852, 2010.
- C. Cartis, N. I. M. Gould, and Ph. L. Toint. Adaptive cubic overestimation methods for unconstrained optimization. Part II: worst-case function-evaluation complexity. *Mathematical Programming, Series A*, **130**(2), 295–319, 2011a.

- C. Cartis, N. I. M. Gould, and Ph. L. Toint. On the evaluation complexity of composite function minimization with applications to nonconvex nonlinear programming. *SIAM Journal on Optimization*, **21**(4), 1721–1739, 2011*b*.
- C. Cartis, N. I. M. Gould, and Ph. L. Toint. An adaptive cubic regularization algorithm for nonconvex optimization with convex constraints and its function-evaluation complexity. *IMA Journal of Numerical Analysis*, **32**(4), 1662–1645, 2012*a*.
- C. Cartis, N. I. M. Gould, and Ph. L. Toint. On the complexity of the steepest-descent with exact linesearches. Technical Report naXys-16-2012, Namur Centre for Complex Systems (naXys), FUNDP-University of Namur, Namur, Belgium, 2012*b*.
- C. Cartis, N. I. M. Gould, and Ph. L. Toint. On the complexity of finding first-order critical points in constrained nonlinear optimization. *Mathematical Programming, Series A*, (to appear), 2013. DOI: 10.1007/s10107-012-0617-9.
- S. Gratton, A. Sartenaer, and Ph. L. Toint. Recursive trust-region methods for multiscale nonlinear optimization. *SIAM Journal on Optimization*, **19**(1), 414–444, 2008.
- L. Grippo, F. Lampariello, and S. Lucidi. A nonmonotone line search technique for Newton’s method. *SIAM Journal on Numerical Analysis*, **23**(4), 707–716, 1986.
- L. Grippo, F. Lampariello, and S. Lucidi. A truncated Newton method with nonmonotone line search for unconstrained optimization. *Journal of Optimization Theory and Applications*, **60**(3), 401–419, 1989.
- Yu. Nesterov. *Introductory Lectures on Convex Optimization*. Applied Optimization. Kluwer Academic Publishers, Dordrecht, The Netherlands, 2004.
- Yu. Nesterov and B. T. Polyak. Cubic regularization of Newton method and its global performance. *Mathematical Programming, Series A*, **108**(1), 177–205, 2006.
- J. Nocedal and S. J. Wright. *Numerical Optimization*. Series in Operations Research. Springer Verlag, Heidelberg, Berlin, New York, 1999.
- J. M. Ortega and W. C. Rheinboldt. *Iterative Solution of Nonlinear Equations in Several Variables*. Academic Press, London, 1970.
- Ph. L. Toint. An assessment of non-monotone linesearch techniques for unconstrained optimization. *SIAM Journal on Scientific and Statistical Computing*, **17**(3), 725–739, 1996.