

Refining asymptotic complexity bounds for nonconvex optimization methods, including why steepest descent is $o(\epsilon^{-2})$ rather than $\mathcal{O}(\epsilon^{-2})$

S. Gratton* and C.-K. Sim[†] and Ph. L. Toint[‡]

16 VIII 2024

Abstract

We revisit the standard “telescoping sum” argument ubiquitous in the final steps of analyzing evaluation complexity of algorithms for smooth nonconvex optimization, and obtain a refined formulation of the resulting bound as a function of the requested accuracy ϵ . While bounds obtained using the standard argument typically are of the form $\mathcal{O}(\epsilon^{-\alpha})$ for some positive α , the refined results are of the form $o(\epsilon^{-\alpha})$. We then explore to which known algorithms our refined bounds are applicable and finally describe an example showing how close the standard and refined bounds can be.

Keywords: Nonlinear optimization, complexity theory, global convergence rates.

1 Introduction

The numerical solution of nonlinear optimization problems often hinges on descent algorithms, that is on algorithms in which a function (the objective function, the residual, the merit function, etc.) is monotonically decreasing in the course of the iterations. The analysis of their iteration (and evaluation) complexity is then typically conducted using a “telescoping sum” argument in which a lower bound of the iteration-wise function decrease is summed in a “telescoping sum” over all iterations. Combining the resulting lower bounds with an upper bound on the total decrease then yields an upper bound on the number of iterations where the iteration-wise decrease is significant, in turn producing the desired upper bound on the algorithm’s worst-case behaviour.

Inspired by an unpublished note [18] of the second author on the steepest descent method, the present paper revisits this telescoping argument, which in turn results in an refined complexity bounds for a large number of known optimization algorithms.

We first describe the refined argument in Section 2 and then investigate to which algorithms the new result is applicable in Section 3. We conclude by presenting, in Section 4, an example showing that the new complexity bounds may be very close to the standard ones.

*Université de Toulouse, INP, IRIT, Toulouse, France. Email: serge.gratton@enseeiht.fr. Work partially supported by 3IA Artificial and Natural Intelligence Toulouse Institute (ANITI), French ”Investing for the Future - PIA3” program under the Grant agreement ANR-19-PI3A-0004”

[†]University of Portsmouth, Hampshire, United Kingdom. Email: chee-khian.sim@port.ac.uk

[‡]NAXYS, University of Namur, Namur, Belgium. Email: philippe.toint@unamur.be

Notation. Given two functions $a(\epsilon)$ and $b(\epsilon)$ with $b(\epsilon) > 0$ depending on a common parameter ϵ tending to zero, we say that $a(\epsilon) = \mathcal{O}(b(\epsilon))$ if and only if there exists a constant $\kappa < +\infty$ such that $\limsup_{\epsilon \rightarrow 0} (a(\epsilon)/b(\epsilon)) \leq \kappa$. We say that $a(\epsilon) = o(b(\epsilon))$ when this limit holds with $\kappa = 0$. If \mathcal{S} is a set, $|\mathcal{S}|$ denotes its cardinality. Finally, $\lambda_{\min}(A)$ denotes the left-most eigenvalue of the symmetric matrix A .

2 A simple result about sequences

In order to discuss our result, we need to consider the situation where a specific optimization algorithm is applied to minimize a smooth possibly nonconvex function f starting from x_0 and producing a sequence of iterates $\{x_k\}$, a sequence of decreasing function values $\{f_k\}$ at these iterates and a sequence of associated optimality measures $\{\omega_k\}$. We also need to consider the set of “successful iterations” $\mathcal{S} = \{k \geq 0 \mid x_{k+1} \neq x_k\}$.

Our results are asymptotic in the sense that we consider these sequences to be infinite and examine how

$$k(\epsilon) = \min\{k \geq 0 \mid \omega_k \leq \epsilon\} \quad (2.1)$$

depends on ϵ when more and more accuracy is requested, that is when ϵ tends to zero. However, since the generation of these sequences will vary across the examples we will consider, we first state our result in a slightly more abstract form.

Theorem 2.1 Let $\{x_k\}$ be a sequence of iterates, $\{f_k = f(x_k)\}$ be a monotonically decreasing sequence bounded below, $\{\omega_k = \omega(x_k)\}$ be a non-negative sequence of optimality measures and let $\mathcal{S} = \{k \geq 0 \mid x_{k+1} \neq x_k\}$ of successful iterations. Suppose also that

$$\mathcal{S} \cap \{k \geq 0 \mid \omega_k = 0\} = \emptyset \quad (2.2)$$

and that

$$f_k - f_{k+1} \geq \kappa_d \omega_k^\beta \quad \text{for } k \in \mathcal{S}, \quad (2.3)$$

where $\kappa_d \in (0, 1]$ and $\beta > 0$ are constants. Suppose also there exist constants $\kappa_a \geq 1$ and $\kappa_b, \kappa_c \geq 0$ such that

$$k \leq \kappa_a |\mathcal{S}_k| + \kappa_b |\log(\omega_k)| + \kappa_c \quad \text{whenever } \omega_k > 0, \quad (2.4)$$

where $\mathcal{S}_k = \mathcal{S} \cap \{0, \dots, k\}$. Then,

$$\lim_{k \rightarrow \infty} \omega_k = 0. \quad (2.5)$$

Moreover, either $k(\epsilon)$ is constant for all ϵ sufficiently small, or

$$k(\epsilon) \leq \kappa_a \max \left[1, \frac{2(f_{\ell(k(\epsilon)-1)} - f_{k(\epsilon)})}{\kappa_d \epsilon^\beta} \right] + \kappa_b |\log(\epsilon)| + \kappa_a + \kappa_c, \quad (2.6)$$

where $\ell(k)$ is the largest index smaller or equal to the median of the indexes in \mathcal{S}_k and

$$\lim_{\epsilon \rightarrow 0} (f_{\ell(k(\epsilon)-1)} - f_{k(\epsilon)-1}) = 0. \quad (2.7)$$

In all cases,

$$k(\epsilon) = o(\epsilon^{-\beta}). \quad (2.8)$$

Proof. Suppose first that there exists a first $k_* \geq 0$ such that $\omega_{k_*} = 0$. Then $k_* \notin \mathcal{S}$. Thus $x_{k+1} = x_k$ and $\omega_{k+1} = 0$, so that, by induction, $\omega_k = 0$ for all $k \geq k_*$, implying (2.5). Moreover, $k(\epsilon) = k_*$ for all $\epsilon < \omega_{k_*-1}$.

Suppose now that $\omega_k > 0$ for all $k \geq 0$. If $|\mathcal{S}|$ is finite, then $\omega_k = \omega_{\min}$ for some $\omega_{\min} > 0$ and all k sufficiently large. As a consequence and given that $|\mathcal{S}_k| \leq |\mathcal{S}| < +\infty$, the right-hand side of (2.4) is bounded. But this is impossible since the left-hand side tends to infinity. Hence $|\mathcal{S}|$ is infinite.

Consider now an arbitrary k for which $|\mathcal{S}_k| \geq 2$. Then $\ell(k)$ is well-defined and tends to infinity with k . We also have, using (2.3) and the definition of $\ell(k)$, that

$$f_{\ell(k)} - f_{k+1} = \sum_{j=\ell(k)}^k (f_j - f_{j+1}) = \sum_{j=\ell(k), j \in \mathcal{S}_k}^k (f_j - f_{j+1}) \geq \frac{1}{2} |\mathcal{S}_k| \kappa_d \min_{j \in \mathcal{S}_k} \omega_j^\beta. \quad (2.9)$$

Moreover, since $\{f_k\}$ is monotonically decreasing and bounded below, it is convergent, and hence

$$\lim_{k \rightarrow \infty} (f_{\ell(k)} - f_{k+1}) = 0. \quad (2.10)$$

Thus the left-hand side of (2.3) tends to zero with k , implying that $\lim_{k \rightarrow \infty, k \in \mathcal{S}} \omega_k = 0$. The definition of $\mathcal{S} \supset \mathcal{S}_k$ then ensures (2.5). As a consequence $k(\epsilon)$ is well-defined for all ϵ sufficiently small. By definition of $k(\epsilon)$ we also know that $\omega_k > \epsilon$ for all $k \leq k(\epsilon) - 1$. Combining this inequality with (2.9), we obtain that

$$|\mathcal{S}_{k(\epsilon)-1}| \leq \frac{2(f_{\ell(k(\epsilon)-1)} - f_{k(\epsilon)})}{\kappa_d \epsilon^\beta}. \quad (2.11)$$

Observe now that $k(\epsilon)$ is non-decreasing when ϵ tends to zero. Given its integer nature, either $k(\epsilon)$ tends to infinity or is constant for all sufficiently small ϵ . In the former case rewriting (2.10) for $k = k(\epsilon) - 1$ gives (2.7) and, because of (2.4) taken at $k(\epsilon)$,

$$k(\epsilon) \leq \kappa_a |\mathcal{S}_{k(\epsilon)}| + \kappa_b \log(\omega_{k(\epsilon)}) + \kappa_c \leq \kappa_a (|\mathcal{S}_{k(\epsilon)-1}| + 1) + \kappa_b |\log(\epsilon)| + \kappa_c, \quad (2.12)$$

which, given (2.11), yields (2.6).

We now prove (2.8). If $k(\epsilon)$ tends to infinity when ϵ tends to zero, (2.8) is obtained by substituting (2.7) in (2.6) and using the fact that $|\log(\epsilon)| = o(\epsilon^{-\beta})$. If $k(\epsilon)$ remains constant, then (2.8) immediately follows from the fact that $\epsilon^{-\beta}$ tends to infinity when ϵ goes to zero. \square

Our assumption (2.2) simply says that, if, luckily, an exact critical point of the desired order is found after finitely many iterations, than the algorithm does not move away. Note that we have chosen $\ell(k)$ above to approximate the index separating \mathcal{S}_k in two parts of same cardinality, but other fixed proportions may of course be used, at the price of modifying the constants in (2.6) and (2.11). Observe also that we could have replaced $|\log(\omega_k)|$ in (2.4) by any positive sequence $\{h(x_k)\}$ such that $h(x_{k(\epsilon)}) = o(\epsilon^{-\beta})$.

3 Application to existing algorithms and associated complexity bounds

We now investigate the consequences of using this simple result in the context where the sequence $\{x_k\}$ is the sequence of iterates generated by specific nonlinear optimization algorithms applied to sufficiently smooth functions that are bounded below. This section only partially explores the resulting refined complexity bounds, focusing on the algorithms described in the comprehensive book [5], but the authors are of course aware that the discussion is incomplete.

3.1 Unconstrained optimization

3.1.1 Steepest descent and other linesearch methods

We start by considering complexity results for linesearch methods for finding first-order critical points, such as those covering steepest descent with Armijo, Goldstein [5, Th. 2.2.2], exact linesearch [5, Th. 2.2.4] or with Nesterov stepsize ([16] and [5, Equation (2.2.5)]). The proof of these results directly involves the “telescoping sum” argument, which we now cast in the context of the previous section by selecting

$$\{f_k\} = \{f(x_k)\}, \quad \omega_k = \|\nabla_x^1 f(x_k)\| \quad \beta = 2, \quad \mathcal{S}_k = \{0, \dots, k\}$$

and κ_d is an algorithm-specific constant proportional to the inverse of the gradient Lipschitz constant. Note that a standard linesearch ensures that $\{f_k\}$ is decreasing in that (2.3) holds at all iterations. Moreover the identity $\mathcal{S}_k = \{0, \dots, k\}$ gives that $\kappa_a = 1$ and $\kappa_b = \kappa_c = 0$.

As a consequence, Theorem 2.1 implies that *the worst-case complexity of all these first-order algorithms (as a function of the accuracy parameter ϵ) is $o(\epsilon^{-2})$ rather than $\mathcal{O}(\epsilon^{-2})$* as stated in the quoted theorems. An illustration for steepest descent is discussed in Section 4.

Interestingly, our technique does not require the complete sequence of function values to satisfy (2.3), but it is enough that these conditions hold, as is the case in the non-monotone “gradient-related” linesearch method discussed in [7], for a subsequence of values at “reference iterations” which is used in the telescoping sum argument. Classical gradient-related linesearch methods [17] are obtained by choosing the memory parameter in this latter method to enforce monotonicity.

3.1.2 Trust-region methods

We may now turn to standard trust-region methods, whose complexity was first considered in [14] and is discussed in [5, Th. 2.3.7 and 3.2.1] (for convergence of first- and second-order methods converging to first-order critical points) and [5, Th. 3.2.6] for convergence to second-order ones. Again the quoted proofs use a “telescoping sum” argument where κ_d an algorithm-specific constant proportional to the inverse of the gradient Lipschitz constant

$$\{f_k\} = \{f(x_k)\}, \quad \omega_k = \|\nabla_x^1 f(x_k)\| \quad \text{or} \quad \omega_k = \max [\|\nabla_x^1 f(x_k)\|, \max(0, \lambda_{\min}(\nabla_x^2 f(x_k)))] ,$$

but we now choose \mathcal{S}_k to be the index set of the “successful iterations”, that is iterations where x_{k+1} differs from x_k and ensuring (2.3). The parameter β now depends on the purpose of the algorithm (finding first- or higher-order critical points) and the degree of the objective’s derivatives used by the algorithm. For standard trust-region methods that seek first-order

critical points, the parameter β is typically equal to two, while it is equal to three if second-order ones are sought. Verifying (2.4) is a little more complicated. [5, Lemma 2.3.1] shows that this inequality holds with ω_k replaced by a lower bound on the trust-region radius. Fortunately, [5, Lem. 2.3.4 and 3.2.5] then state that this lower bound is itself bounded below by ω_k or ϵ (for $k = k(\epsilon)$), hence providing the desired inequality.

We may thus again apply our results to revisit all these proofs. For the search of first-order critical points, this gives $o(\epsilon^{-2})$ rather than $\mathcal{O}(\epsilon^{-2})$ complexity bounds as a function of ϵ . The bounds for finding second-order points are similarly refined to $o(\epsilon^{-3})$ rather than $\mathcal{O}(\epsilon^{-3})$. In the same vein, we may even consider trust-region methods for delivering critical points of order higher than two [5, Th. 12.2.5] and obtain $o(\epsilon^{-(q+1)})$ worst-case complexity to compute q -th order critical points. Finally, the global rates of convergence of TRqIDA and TRqEDA trust-region variants for noisy problems may also be refined in the same way (see [5, Th. 13.1.8, 13.3.4] together with [5, Lem. 13.1.1 and 13.1.4]).

But we may also consider more elaborate trust-region-like algorithms, such as TR ϵ [8] (whose complexity proofs can be found in [5, Th. 3.4.5 and 3.4.6]), TRACE [9] (see [5, Th. 3.4.11 and 3.4.12] for proofs) or the Birgin-Martinez proposal [1]. These methods achieve a complexity bound using $\beta = 3/2$ when first-order points are sought. Note that a specific result [5, Lem. 3.4.8 and 3.4.10] is needed for the second of these methods to yield (2.4). Since these methods have a better ϵ -order complexity, we now deduce that it is now $o(\epsilon^{-3/2})$ rather than $\mathcal{O}(\epsilon^{-3/2})$ for finding first-order critical points, and $o(\epsilon^{-3})$ rather than $\mathcal{O}(\epsilon^{-3})$ to find second-order ones.

3.2 Adaptive regularization methods

The case of adaptive regularization methods is quite similar to that of trust-region algorithms. Again

$$\{f_k\} = \{f(x_k)\}, \quad \omega_k = \text{a } q\text{-th order criticality measure,}$$

κ_d is an algorithm-specific constant and \mathcal{S}_k is the index set of the “successful iterations”. The bound (2.4) is now guaranteed by [5, Lem 2.4.1 and 2.4.2] with $\kappa_b = 0$ and β again depends on which type of critical points are sought and the degree of derivatives used. Because a specific discussion of every case may quickly become cumbersome, we only list, in Table 1, the algorithms of interest, pointers to the relevant proofs, criticality order q and associated refined complexity bounds resulting from Theorem 2.1.

The proofs for the AN2C algorithms [12], using an alternative regularization of Newton’s method, are more involved because \mathcal{S}_k is then the union of smaller sets, but again rely on “telescoping sums” for subsets of iterations, [12, Lem. 1 and 4] being used to ensure (2.4) in this case. The AR1pGN and AR2GN algorithms proposed in [15] allows the use of a general nonsmooth regularization, and (2.4) is ensured by [15, Lem. 2.4 and 3.3] in this case.

As it turns out, the proofs listed in Table 1 are themselves templates for the complexity proofs of variants of the adaptive regularization that exploit problem structure. Again discussing every case would be too cumbersome, but we refer the reader to [6] for a specialized algorithm for least Euclidean distance optimization, to [5, Th. 14.1.10] for a variant designed for the minimization of possibly non-smooth composite objectives, to [5, Th. 13.1.19 and 13.3.8] (together with [5, Lem 13.1.9]) for noise-tolerant variants or to [4] for an algorithm exploiting finite-differences approximations to derivatives, including the derivative-free case.

Algo.	Proof	Critic. order	Refined complexity
AR1	[5, Th. 2.4.3]	1rst	$o(\epsilon^{-2})$
AR2	[5, Th. 3.3.4]	1rst	$o(\epsilon^{-3/2})$
AR2	[5, Th. 3.3.9]	2nd	$o(\epsilon^{-3})$
AR p	[5, Th. 4.1.5]	1rst, 2nd, 3rd	$o(\epsilon^{-(p+1)/(p+1-q)})$
AR qp	[5, Th. 12.2.14]	1rst, 2nd	$o(\epsilon^{-(p+1)/(p-q+1)})$
AR qp	[5, Th. 12.2.14]	q -th, $q > 2$	$o(\epsilon^{-q(p+1)/p})$
AR qp IDA	[5, Th. 13.1.19]	1rst, 2nd	$o(\epsilon^{-(p+1)/(p-q+1)})$
AR qp IDA	[5, Th. 13.1.19]	q -th, $q > 2$	$o(\epsilon^{-q(p+1)/p})$
AR qp EDA	[5, Th. 13.3.8]	1rst, 2nd	$o(\epsilon^{-(p+1)/(p-q+1)})$
AR qp EDA	[5, Th. 13.3.8]	q -th, $q > 2$	$o(\epsilon^{-q(p+1)/p})$
AN2C	[12, Th. 1]	1rst	$o(\epsilon^{-3/2})$
AN2C	[12, Th. 2]	2nd	$o(\epsilon^{-3})$
AR1 p GN	[15, Th. 3.5]	1rst	$o(\epsilon^{-(p+1)/p})$
AR2GN	[15, Th. 4.5]	1rst, 2nd	$o(\epsilon^{-(p+1)/(p+1-q)})$

Table 1: Refined complexity bound for unconstrained adaptive regularization algorithms

3.3 Direct-search methods

Finally, direct-search methods for minimization may also be considered. In [19, Th. 3.1] the telescoping sum argument is again explicitly used to prove a worst-case complexity bound for this important class of derivative-free methods. In this case, \mathcal{S} is the set of successful iterations (as for trust-region and adaptive regularization algorithms), $\{f_k\} = \{f(x_k)\}$, $\omega_k = \|\nabla_x^1 f(x_k)\|$, $\beta = 2$, κ_d a constant involving the square of the gradient's Lipschitz constant. The bound (2.4) is obtained by [19, Lem. 3.2] for $k = k(\epsilon)$ (as needed to derive (2.12)). The complexity bound in $\mathcal{O}(\epsilon^{-2})$ of [19, Cor. 3.1] can then be refined to $o(\epsilon^{-2})$.

3.4 Algorithms for constrained problems

Because methods for unconstrained optimization do occur as crucial ingredients of several algorithms for the constrained case, the refined complexity bounds for the former may translate into refined complexity bounds for the latter. The easiest situation is when considering ‘‘simple’’ constraints, i.e. when the constraints define a convex feasible set onto which projection is computationally affordable (including, for example, the ubiquitous problem of minimizing a function subject to simple bounds on the variables). In this case, the evaluation complexity bounds for unconstrained problems are often unmodified (when considering their order as a function of ϵ) compared with their unconstrained counter-parts, and the techniques of proof to establish them are directly derived from the unconstrained setting, except for the use of criticality measures that are suitable for constrained problems. See for instance [5, Th. 6.2.3], where the $\mathcal{O}(\cdot)$ bounds for first-, second- and third order critical points may now be refined to $o(\cdot)$.

Finally, unconstrained or bound-constrained methods and the techniques to prove their complexity are often instrumental in the analysis of algorithms for more general nonlinear constraints (for instance for ‘‘restoration’’ or ‘‘feasibility’’ phases, where one minimizes the violation of the nonlinear constraints, essentially using algorithms for unconstrained problems).

Resulting bounds for the whole constrained algorithm may then be refined along the lines described above (see, for instance, [5, Th. 7.2.2 and 7.2.6] leading to [5, Th. 7.2.7]).

4 How close are the refined and standard bounds?

Having discussed refined bounds for a significant selection of algorithms, we now take a step back and investigate how much the refined and standard bounds differ by looking at a particular example. This example is univariate and built along the lines of [5, Th. 2.2.3] for steepest descent. Sequences of iterates $\{x_k\}$, function values $\{f_k\}$, gradient values $\{g_k\}$ and steps $\{s_k\}$ are first constructed to illustrate the bound, and standard Hermite theory is then invoked to show the existence of a suitable function interpolating these values.

Define, for $k \geq 0$ and some fixed constant $\delta > 0$,

$$g_0 = -2 \quad \text{and} \quad g_k = -\frac{1}{k^{\frac{1}{2}+\delta}} \quad (k > 0), \quad (4.1)$$

$$f_0 = \zeta(1+2\delta) > 1, \quad f_1 = f_0 - 4\alpha \quad \text{and} \quad f_{k+1} = f_k - \frac{\alpha}{k^{1+2\delta}} \quad (k > 0) \quad (4.2)$$

and

$$x_0 = 0, \quad x_1 = 2\alpha \quad \text{and} \quad x_{k+1} = x_k + \frac{\alpha}{k^{\frac{1}{2}+\delta}} \quad (k > 0). \quad (4.3)$$

for some $\alpha \in (0, 1]$, where $\zeta(\cdot)$ is the Riemann function. By definition of this function, we then have that $f_k \geq 0$ for all $k \geq 0$, so the sequence $\{f_k\}$ is strictly decreasing and bounded below, and hence convergent to some limit value $f_{\text{lim}} \geq 0$. As a consequence we have that

$$\lim_{k \rightarrow \infty} (f_{\lfloor k/2 \rfloor} - f_k) = 0. \quad (4.4)$$

Moreover, let

$$s_0 = 2\alpha \quad \text{and} \quad s_k = \frac{\alpha}{k^{\frac{1}{2}+\delta}}. \quad (4.5)$$

A simple calculation then shows that $k(\epsilon)$ as defined by (2.1) satisfies

$$k(\epsilon) = \lceil k_\epsilon \rceil \quad \text{where} \quad k_\epsilon = \epsilon^{-\frac{1}{\frac{1}{2}+\delta}} = \epsilon^{-2} \epsilon^{\frac{4\delta}{1+2\delta}} = o(\epsilon^{-2}) \quad (4.6)$$

and thus

$$k(\epsilon) \leq k_\epsilon + 1 = o(\epsilon^{-2}), \quad (4.7)$$

which is (2.8) (note that $k(\epsilon)$ tends to infinity when ϵ tends to zero because of (4.1)). But (4.2), (4.6) and (4.7) together give that

$$f_{\lfloor (k(\epsilon)-1)/2 \rfloor} - f_{k(\epsilon)} \geq \sum_{k=\lfloor (k(\epsilon)-1)/2 \rfloor}^{k(\epsilon)-1} \frac{\alpha}{k^{1+2\delta}} \geq \frac{1}{2} |\mathcal{S}_{k(\epsilon)-1}| \frac{\alpha}{k_\epsilon^{1+2\delta}} = \frac{1}{2} |\mathcal{S}_{k(\epsilon)-1}| \frac{\alpha}{\epsilon^{-2}}$$

and thus the stronger bound (2.11) also holds. Now, because

$$|f_{k+1} - f_k - g_k s_k| = 0 \leq s_k^2$$

and

$$|g_0 - g_1| = |2 - 1| \leq \frac{1}{\alpha} s_0 \quad |g_{k+1} - g_k| = \frac{1}{k^{\frac{1}{2}+\delta}} - \frac{1}{(k+1)^{\frac{1}{2}+\delta}} \leq \frac{1}{\alpha} s_k,$$

we may apply Hermite's theorem [5, Th. A.9.1] with $\kappa_f = f_0$ on each interval $[x_k, x_{k+1}]$, defining a cubic polynomial interpolating f_k, f_{k+1}, g_k and g_{k+1} . Combining these polynomials for successive intervals, we obtain a cubic piecewise polynomial f which is continuously differentiable from $[0, +\infty)$ into \mathbb{R} and whose gradient is Lipschitz continuous with a constant L only depending on $\kappa_f > 1$ and $1/\alpha$. In addition, for all $k \geq 0$,

$$f(x_k) = f_k \text{ and } \nabla_x f(x_k) = g_k$$

and $f(x)$ is bounded below on $[0, +\infty)$. It is then easy to extend this function on the left of the origin without altering these properties by defining $f(x) = f_0 - 2x$ for $x < 0$. The left panel of Figure 1 shows the (deceptively innocuous looking and barely nonconvex) graph of f in the interval $[x_0, x_4]$, where $\alpha = 0.1$, $\delta = 0.001$ and $\epsilon = 0.01$ (note that $\zeta(1.002) \approx 500.577$). The middle panel shows the graph of its (continuous but non-monotone) gradient and the right one its (discontinuous but bounded) Hessian.

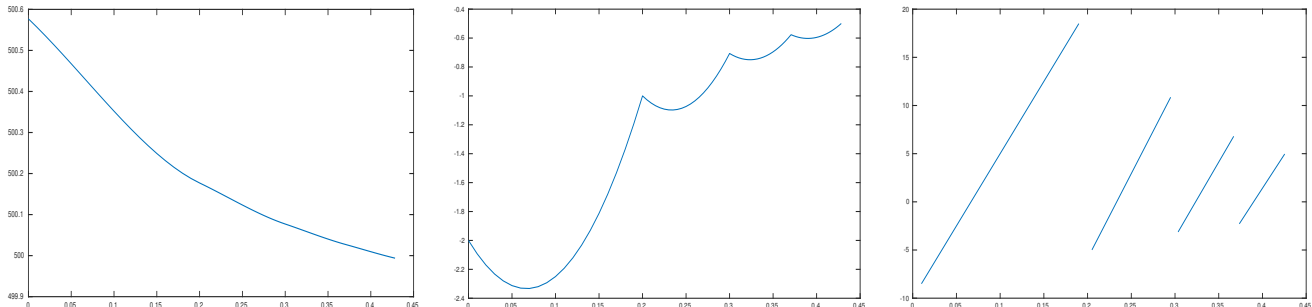


Figure 1: The functions f (left), $\nabla_x^1 f$ (middle) and $\nabla_x^2 f$ (right) in $[x_0, x_4]$

As a consequence of the above argument, we may interpret the sequences $\{x_k\}$, $\{f_k\}$ and $\{g_k\}$ as the result of a steepest descent algorithm using an Armijo linesearch with initial stepsize α , applied to the univariate function f and starting from x_0 . The initial stepsize is always acceptable for the linesearch because sufficient decrease (αg_k^2) is achieved for the initial stepsize at every iteration. In view of (4.7), we have thus verified that, as expected, our complexity bound in $o(\epsilon^{-2})$ holds for the steepest descent algorithm.

Now looking at (4.7), we also see that *this bound can be arbitrarily close to standard worst-case bound* in $\mathcal{O}(\epsilon^{-2})$ when δ is close to zero. A similar conclusion also holds for the other cases discussed in Section 3.

5 Conclusion

We have revisited the last step of the worst-case complexity proofs for nonlinear optimization algorithms and obtained refined theoretical bounds. We have then considered a few of the many cases where these proofs can be refined, but the idea can clearly be applied more widely. We have also shown that, although better, the refined bound may be arbitrarily close to the standard one.

We note that our asymptotic results do not contradict the non-asymptotic lower complexity bounds proved in [3] and [5, Th. 2.2.3, 2.2.16 and 12.2.17]. Indeed these latter bounds depend on examples where a function is constructed such that convergence of the relevant algorithm is exactly as slow as specified by the standard $\mathcal{O}(\cdot)$ bound. However, these functions

explicitly depend on ϵ , which prevents taking the limit for ϵ tending to zero, as we have done above. The situation is similar for the example proposed in [2], where the function on which slow convergence occurs is defined in a space whose dimension depends on ϵ .

Of course, not all convergence proofs (and algorithms) are concerned. Notable exceptions include complexity proofs for measure-dominated problems (see [5, Section 5.3]) because proofs in this context do not directly rely on the telescoping sum argument. The case of objective-function free (OFFO) algorithms, among which many stochastic methods (see, for example, [13, 20, 10, 11]), is less clear because the relevant complexity proofs typically involve telescoping sums along with other potentially dominating terms.

While the refined bounds are interesting, they remain of a fairly generic nature, as we have verified in Section 3. It remains an open question whether they can be refined further (maybe by quantifying the numerator of the right-hand side of (2.11)) for specific methods.

Acknowledgements

The authors wish to thank Oliver Hinder for an interesting discussion at ISMP 2024. The third author also gratefully acknowledges the friendly partial support of ANTI (Toulouse).

References

- [1] E. G. Birgin and J. M. Martínez. On regularization and active-set methods with complexity for constrained optimization. *SIAM Journal on Optimization*, 28(2):1367–1395, 2028.
- [2] Y. Carmon, J. C. Duchi, O. Hinder, and A. Sidford. Lower bounds for finding stationary points I. *Mathematical Programming, Series A*, 184:71–120, 2020.
- [3] 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.
- [4] C. Cartis, N. I. M. Gould, and Ph. L. Toint. On the oracle complexity of first-order and derivative-free algorithms for smooth nonconvex minimization. *SIAM Journal on Optimization*, 22(1):66–86, 2012.
- [5] C. Cartis, N. I. M. Gould, and Ph. L. Toint. *Evaluation complexity of algorithms for nonconvex optimization*. Number 30 in MOS-SIAM Series on Optimization. SIAM, Philadelphia, USA, June 2022.
- [6] C. Cartis, N. I. M. Gould, and Ph. L. Toint. Strong evaluation complexity bounds for arbitrary-order optimization of nonconvex nonsmooth composite functions. In *Invited Lectures, Proceedings of the 2022 International Conference of Mathematicians (ICM 2022), St Petersburg*, pages 5266–5289. European Mathematical Society (EMS), 2022.
- [7] C. Cartis, Ph. R. Sampaio, and Ph. L. Toint. Worst-case complexity of first-order non-monotone gradient-related algorithms for unconstrained optimization. *Optimization*, 64(5):1349–1361, 2015.
- [8] F. Curtis, D. Robinson, C. Royer, and S. J. Wright. Trust-region Newton-CG with strong second-order complexity guarantees for nonconvex optimization. *SIAM Journal on Optimization*, 31(1):518–544, 2021.
- [9] F. E. Curtis, D. P. Robinson, and M. Samadi. An inexact regularized Newton framework with a worst-case iteration complexity of $O(\epsilon^{-3/2})$ for nonconvex optimization. *IMA Journal of Numerical Analysis*, 00:1–32, 2018.
- [10] A. Défossez, L. Bottou, F. Bach, and N. Usunier. A simple convergence proof for Adam and Adagrad. *Transactions on Machine Learning Research*, October 2022.
- [11] S. Gratton, S. Jerad, and Ph. L. Toint. A stochastic objective-function-free adaptive regularization method with optimal complexity. arXiv:2407.08018, 2024.
- [12] S. Gratton, S. Jerad, and Ph. L. Toint. Yet another fast variant of Newton’s method for nonconvex optimization. *IMA Journal of Numerical Analysis*, (to appear), 2024.
- [13] S. Gratton, C. W. Royer, L. N. Vicente, and Z. Zhang. Direct search based on probabilistic descent. *SIAM Journal on Optimization*, 25(3):1515–1541, 2015.

- [14] 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.
- [15] S. Gratton and Ph. L. Toint. Adaptive regularization minimization algorithms with non-smooth norms. *IMA Journal of Numerical Analysis*, 43(2):920–949, 2023.
- [16] Yu. Nesterov. *Introductory Lectures on Convex Optimization*. Applied Optimization. Kluwer Academic Publishers, Dordrecht, The Netherlands, 2004.
- [17] G. A. Shultz, R. B. Schnabel, and R. H. Byrd. A family of trust-region-based algorithms for unconstrained minimization with strong global convergence properties. *SIAM Journal on Numerical Analysis*, 22(1):47–67, 1985.
- [18] C.-K. Sim. A $o(1/\rho^2)$ -first order algorithm on composite nonconvex optimization problems. <https://drive.google.com/file/d/1NVsQcrW0TL3x0mbPeP-edFeJi6Yfwlv/view?pli=1>, 2021.
- [19] L. N. Vicente. Worst case complexity of direct search. *EURO Journal on Computational Optimization*, 1:143–153, 2013.
- [20] R. Ward, X. Wu, and L. Bottou. Adagrad stepsizes: sharp convergence over nonconvex landscapes. In *Proceedings in the International Conference on Machine Learning (ICML2019)*, 2019.