

Quadratic regularization with cubic descent for unconstrained optimization*

E. G. Birgin[†] J. M. Martínez[‡]

October 27, 2016

Abstract

Cubic-regularization and trust-region methods with worst-case first-order complexity $O(\varepsilon^{-3/2})$ and worst-case second-order complexity $O(\varepsilon^{-3})$ have been developed in the last few years. In this paper it is proved that the same complexities are achieved by means of a quadratic-regularization method with a cubic sufficient-descent condition instead of the more usual predicted-reduction based descent. Asymptotic convergence and order of convergence results are also presented. Finally, some numerical experiments comparing the new algorithm with a well-established quadratic regularization method are shown.

Key words: Nonlinear programming, unconstrained minimization, quadratic regularization, cubic descent, complexity.

1 Introduction

Assume that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is possibly nonconvex and smooth for all $x \in \mathbb{R}^n$. We will consider the unconstrained minimization problem given by

$$\text{Minimize } f(x). \tag{1}$$

In the last decade, many works has been devoted to analyze iterative algorithms for solving (1) from the point of view of their time complexity. See, for example, [2, 4, 5, 6, 8, 11, 14, 19, 20]. A review of complexity results for the convex case, in addition to novel techniques, can be found in [12].

Given arbitrary tolerances $\varepsilon_g > 0$ and $\varepsilon_H > 0$ the question is about the amount of iterations, functional and derivative evaluations, that are necessary to achieve an approximate solution defined by $\|\nabla f(x)\| \leq \varepsilon_g$ or by $\|\nabla f(x)\| \leq \varepsilon_g$ plus $\lambda_1(\nabla^2 f(x)) \geq -\varepsilon_H$, where $\lambda_1(\nabla^2 f(x))$ represents the smallest eigenvalue of $\nabla^2 f(x)$.

*This work has been partially supported by FAPESP (grants 2010/10133-0, 2013/03447-6, 2013/05475-7, 2013/07375-0, and 2014/18711-3) and CNPq (grants 309517/2014-1 and 303750/2014-6).

[†]Department of Computer Science, Institute of Mathematics and Statistics, University of São Paulo, Rua do Matão, 1010, Cidade Universitária, 05508-090, São Paulo, SP, Brazil. e-mail: egbirgin@ime.usp.br

[‡]Department of Applied Mathematics, Institute of Mathematics, Statistics, and Scientific Computing (IMECC), State University of Campinas, 13083-859 Campinas SP, Brazil. e-mail: martinez@ime.unicamp.br

In general, gradient-based methods exhibit complexity $O(\varepsilon_g^{-2})$ [4], which means that there exists a constant c that only depends on the characteristics of the problem, algorithmic parameters, and, of course, the initial approximation, such that the effort required to achieve $\|\nabla f(x)\| \leq \varepsilon_g$ for a bounded-below objective function f is at most c/ε_g^2 . This bound is sharp for all gradient-based methods [4]. Complexity results for quasi-Newton methods are available in [14]. Surprisingly, Newton’s method with the classical trust-region strategy does not exhibit better complexity than $O(\varepsilon_g^{-2})$ either [4]. The same example used in [4] to prove this fact can be applied to Newton’s method with standard quadratic regularization. On the other hand, Newton’s method employing cubic regularization [15] for obtaining sufficient descent at each iteration exhibits the better complexity $O(\varepsilon_g^{-3/2})$ (see [5, 6, 19, 20]).

The best known practical algorithm for unconstrained optimization with worst-case evaluation complexity $O(\varepsilon_g^{-3/2})$ to achieve first-order stationarity and complexity $O(\varepsilon_g^{-3/2} + \varepsilon_H^{-3})$ to achieve second-order stationarity, defined by Cartis, Gould, and Toint in [5] and [6], uses cubic regularization and a descent criterion based on the comparison of the actual reduction of the objective function and the reduction predicted by a quadratic model. A non-standard trust-region method with the same complexity properties due to Curtis, Robinson, and Samadi [8] employs a cubic descent criterion for accepting trial increments. In [2], the essential ideas of ARC [5, 6] were extended in order to introduce high-order methods in which a p -th Taylor approximation ($p \geq 2$) plus a $(p + 1)$ -th regularization term is minimized at each iteration. In these methods, $O(\varepsilon_g^{-(p+1)/p})$ evaluation complexity for first-order stationarity is obtained also using the actual-versus-predicted-reduction descent criterion. However, it is rather straightforward to show that this criterion can be replaced by a $(p + 1)$ -th descent criterion (i.e. $f(x^{k+1}) \leq f(x^k) - \alpha\|x^{k+1} - x^k\|^{p+1}$) in order to obtain the same complexity results. Moreover, the $(p + 1)$ -th descent criterion (cubic descent in the case $p = 2$) seems to be more naturally connected with the Taylor approximation properties that are used to prove complexity. Cubic descent was also used in [19] in a variable metric method that seeks to achieve good practical global convergence behavior.

In the trust-region example exhibited in [4], the unitary Newtonian step is accepted at every iteration since it satisfies the adopted sufficient descent criterion. This criterion requires that the function descent (actual reduction) should be better than a fraction of the predicted descent provided by the quadratic model (predicted reduction). However, if, instead of this condition, one requires functional descent proportional to $\|s\|^3$, where s is the increment given by the model minimization, the given example does not stand anymore. This state of facts led us to the following theoretical question: Would it be possible to obtain worst-case evaluation complexities $O(\varepsilon_g^{-3/2})$ and $O(\varepsilon_g^{-3/2} + \varepsilon_H^{-3})$ using cubic descent to accept trial increments but only quadratic regularization in the subproblems? It is well known that the set of solutions of the cubic regularization problem, for different values of the cubic regularization parameter, coincides with the set of solutions of the quadratic regularization problem. This set of solutions is known as Levenberg-Marquardt path. Then, the question is whether it is possible to visit different points of the Levenberg-Marquardt path, increasing the quadratic regularization parameter (or maintaining its value as in the “hard case”) in such a way that the main arguments for obtaining the desired complexity results are maintained. Looking for the answer to this question, we arrived to a sophisticated procedure for traveling along the Levenberg-Marquardt path that satisfies the

desired complexity requirements. Basically, the so called “hard case” must be carefully taken into account and “the first” regularization parameter tried at each iteration should be chosen between suitable safety bounds. Only for choosing “the first” regularization parameter it could be necessary to solve more than one linear system per function evaluation. As a result, we came up with a geometrically-inspired one-page easy-to-implement method that can be described in no more than six simple steps.

Although, in principle, we did not aim to produce a practical algorithm, we adopted the point of view that algorithmic features that are responsible for theoretical improvements frequently lead also to practical advantages. Therefore, we implemented the new method and compared it against a well established quadratic regularization method for unconstrained optimization introduced in [16]. The results seem to be interesting for future research.

The rest of this paper is organized as follows. A model algorithm with cubic descent is described in Section 2. An implementable version of the algorithm is introduced in Section 3. Well-definiteness and complexity results are presented in Sections 4 and 5, respectively. Local convergence results are given in Section 6. Numerical experiments are presented in Section 7; while final remarks are given in the last section.

Notation. The symbol $\|\cdot\|$ denotes the Euclidean norm of vectors and the subordinate matricial norm. We denote $g(x) = \nabla f(x)$, $H(x) = \nabla^2 f(x)$, and, sometimes, $g^k = g(x^k)$ and $H_k = H(x^k)$. If $a \in \mathbb{R}$, $[a]_+ = \max\{a, 0\}$. If $a_1, \dots, a_n \in \mathbb{R}$, $\text{diag}(a_1, \dots, a_n)$ denotes the $n \times n$ diagonal matrix whose diagonal entries are a_1, \dots, a_n . The notation $[x]_j$ denotes the j th component of a vector x whenever the simpler notation x_j might lead to confusion.

2 Model algorithm

The following algorithm establishes a general framework for minimization schemes that use cubic descent. At each iteration k , we compute an increment s^k such that $f(x^k + s^k) \leq f(x^k) - \alpha \|s^k\|^3$. In principle, this is not very useful because even $s^k = 0$ satisfies this descent condition. However, in Theorem 2.1, we show that under the additional condition (3), the algorithm satisfies suitable stopping criteria. As a consequence, practical algorithms should aim to achieve (2) and (3) simultaneously.

Algorithm 2.1

Let $x^0 \in \mathbb{R}^n$ and $\alpha > 0$ be given. Initialize $k \leftarrow 0$.

Step 1. Compute s^k such that

$$f(x^k + s^k) \leq f(x^k) - \alpha \|s^k\|^3. \quad (2)$$

Step 2. Define $x^{k+1} = x^k + s^k$, set $k \leftarrow k + 1$, and go to Step 1.

The theorems below establish that, under suitable assumptions, every limit point of the sequence generated by Algorithm 2.1 is second-order stationary and provide an upper bound

on the number of iterations that Algorithm 2.1 requires to achieve a target objective functional value or to find an approximate first- or second-order stationary point.

Lemma 2.1 *Assume that the objective function f is twice continuously differentiable and that there exist $\gamma_g > 0$ and $\gamma_H > 0$ such that, for all $k \in \mathbb{N}$, the increment s^k computed at Step 1 of Algorithm 2.1 satisfies*

$$\sqrt{\frac{\|g^{k+1}\|}{\gamma_g}} \leq \|s^k\| \text{ and } \frac{[-\lambda_{1,k}]_+}{\gamma_H} \leq \|s^k\|, \quad (3)$$

where $\lambda_{1,k}$ stands for the smallest eigenvalue of H_k . Then, it follows that

$$f(x^{k+1}) \leq f(x^k) - \max \left\{ \left(\frac{\alpha}{\gamma_g^{3/2}} \right) \|g^{k+1}\|^{3/2}, \left(\frac{\alpha}{\gamma_H^3} \right) [-\lambda_{1,k}]_+^3 \right\}.$$

Proof: The result follows trivially from (2), (3), and the fact that, at Step 2 of Algorithm 2.1, x^{k+1} is defined as $x^{k+1} = x^k + s^k$. \square

Theorem 2.1 *Let $f_{\min} \in \mathbb{R}$, $\varepsilon_g > 0$, and $\varepsilon_H > 0$ be given constants, assume that the hypothesis of Lemma 2.1 hold, and let $\{x^k\}_{k=0}^\infty$ be the sequence generated by Algorithm 2.1. Then, the cardinality of the set of indices*

$$K_g = \left\{ k \in \mathbb{N} \mid f(x^k) > f_{\min} \text{ and } \|g^{k+1}\| > \varepsilon_g \right\} \quad (4)$$

is, at most,

$$\left\lfloor \frac{1}{\alpha} \left(\frac{f(x^0) - f_{\min}}{(\varepsilon_g/\gamma_g)^{3/2}} \right) \right\rfloor; \quad (5)$$

while the cardinality of the set of indices

$$K_H = \left\{ k \in \mathbb{N} \mid f(x^k) > f_{\min} \text{ and } \lambda_{1,k} < -\varepsilon_H \right\} \quad (6)$$

is, at most,

$$\left\lfloor \frac{1}{\alpha} \left(\frac{f(x^0) - f_{\min}}{(\varepsilon_H/\gamma_H)^3} \right) \right\rfloor. \quad (7)$$

Proof: From Lemma 2.1, it follows that at every time an iterate x^k is such that $\|g^{k+1}\| > \varepsilon_g$ the value of f decreases at least $\alpha(\varepsilon_g/\gamma_g)^{3/2}$; while at every time an iterate x^k is such that $\lambda_{1,k} < -\varepsilon_H$ the value of f decrease at least $\alpha(\varepsilon_H/\gamma_H)^3$. Therefore, the thesis follows from the fact that, by (2), $\{f(x^k)\}_{k=0}^\infty$ is a non-increasing sequence. \square

Corollary 2.1 *Let $f_{\min} \in \mathbb{R}$, $\varepsilon_g > 0$, and $\varepsilon_H > 0$ be given constants and assume that the hypothesis of Lemma 2.1 hold. Algorithm 2.1 requires $O(\varepsilon_g^{-3/2})$ iterations to compute x^k such that*

$$f(x^k) \leq f_{\min} \text{ or } \|g^{k+1}\| \leq \varepsilon_g;$$

it requires $O(\varepsilon_{\text{H}}^{-3})$ iterations to compute x^k such that

$$f(x^k) \leq f_{\min} \text{ or } \lambda_{1,k} \geq -\varepsilon_{\text{H}};$$

and it requires $O(\varepsilon_g^{-3/2} + \varepsilon_{\text{H}}^{-3})$ iterations to compute x^k such that

$$f(x^k) \leq f_{\min} \text{ or } \left(\|g^{k+1}\| \leq \varepsilon_g \text{ and } \lambda_{1,k} \geq -\varepsilon_{\text{H}} \right).$$

Corollary 2.2 *Assume that the hypothesis of Lemma 2.1 hold and let $\{x^k\}_{k=0}^{\infty}$ be the sequence generated by Algorithm 2.1. Then, if the objective function f is bounded below, we have that*

$$\lim_{k \rightarrow \infty} \|g(x^k)\| = 0 \text{ and } \lim_{k \rightarrow \infty} [-\lambda_{1,k}]_+ = 0.$$

Proof: Assume that $\lim_{k \rightarrow \infty} \|g(x^k)\| \neq 0$. This means that there exists $\varepsilon > 0$ and \mathbb{K} , an infinite subsequence of \mathbb{N} , such that $\|g^k\| > \varepsilon$ for all $k \in \mathbb{K}$. Since f is bounded below, this contradicts Theorem 2.1. The second part is analogous. \square

Corollary 2.3 *Assume that the hypothesis of Lemma 2.1 hold. Then, if the objective function f is bounded below, every limit point x^* of the sequence $\{x^k\}_{k=0}^{\infty}$ generated by Algorithm 2.1 is such that $\nabla f(x^*) = 0$ and $\nabla^2 f(x^*)$ is positive semidefinite.*

Proof: This corollary follows from Corollary 2.2 by continuity of ∇f and $\nabla^2 f$. \square

3 Implementable algorithm

Algorithm 2.1 presented in the previous section is a “model algorithm” in the sense that it does not prescribe a way to compute the step s^k satisfying (2). This will be the subject of the present section. Algorithm 3.1 is almost identical to Algorithm 2.1 with the sole difference that it uses Algorithm 3.2 to compute s^k . Lemma 4.1 shows that Algorithm 3.2 is well defined and Lemma 4.4 shows that the step s^k computed by Algorithm 3.2 satisfies the hypothesis (3) of Lemma 2.1. In the following section, it will be shown that Algorithm 3.2 computes s^k using $O(1)$ evaluations of f (and a single evaluation of g and H at the current iterate x^k). This implies that the complexity results on the number of iterations of the model Algorithm 2.1 also apply to the number of iterations and evaluations of f and its first- and second-order derivatives performed by Algorithm 3.1–3.2.

Algorithm 3.1

Let $x^0 \in \mathbb{R}^n$, $\alpha > 0$, and $M > 0$ be given. Initialize $k \leftarrow 0$.

Step 1. Use Algorithm 3.2 to compute $s \in \mathbb{R}^n$ satisfying

$$f(x^k + s) \leq f(x^k) - \alpha \|s\|^3 \tag{8}$$

and define $s^k = s$.

Step 2. Define $x^{k+1} = x^k + s^k$, set $k \leftarrow k + 1$, and go to Step 1.

Algorithm 3.2 below describes the way in which the increment s^k is computed. For that purpose, different trial increments are tried along the set of solutions $s(\mu)$ of

$$\text{Minimize } \langle g^k, s \rangle + \frac{1}{2} s^T (H_k + [-\lambda_{1,k}]_+) s + \frac{\mu}{2} \|s\|^2, \quad (9)$$

where $\lambda_{1,k}$ is the lowest eigenvalue of H_k , for different values of the regularizing parameter $\mu \geq 0$. Roughly speaking, μ will be increased in a controlled way up to the verification of the sufficient descent condition (8) with $s = s(\mu)$.

The geometry of the set of solutions of (9), many times called Levenberg-Marquardt path, will be fully exploited. When μ tends to infinity, $s(\mu)$ tends to 0 describing a curve tangent to $-g^k$. If H_k is positive definite then the Levenberg-Marquardt path is a bounded curve that joins $s = 0$ with the Newtonian step $s = -H_k^{-1}g^k$. If H_k is not positive definite and the system $[H_k + [-\lambda_{1,k}]_+I]s = -g^k$ is not compatible then the Levenberg-Marquardt path is an unbounded curve that, as μ tends to 0, becomes tangent to an affine subspace generated by an eigenvector of H_k associated with $\lambda_{1,k}$ (see Figure 1a). The case in which H_k is not positive definite but $[H_k + [-\lambda_{1,k}]_+I]s = -g^k$ is compatible is called ‘‘hard case’’ in the trust-region literature [7]. In the hard case, the Levenberg-Marquardt path is constituted by two parts. The first part, that corresponds to $\mu > 0$, is a bounded curve that joins the minimum-norm solution of $[H_k + [-\lambda_{1,k}]_+I]s = -g^k$ with $s = 0$. The second part, that corresponds to $\mu = 0$, is given by the infinitely many solutions to the system $[H_k + [-\lambda_{1,k}]_+I]s = -g^k$ (note that $s(\mu)$ is *not* univocally determined in this case as it is in the other cases). This set of infinitely many solutions form an affine subspace that contains $-[H_k + [-\lambda_{1,k}]_+I]^\dagger g^k$ and is spanned by the eigenvectors of H_k associated with $\lambda_{1,k}$ (see Figure 1b). Usually, one restricts this affine subspace to the line $-[H_k + [-\lambda_{1,k}]_+I]^\dagger g^k + tv$ with $t \in \mathbb{R}$, where v is one of the eigenvectors associated with $\lambda_{1,k}$. According to this geometry, in each situation we will define appropriate strategies to travel along the Levenberg-Marquardt path.

At a first glance, Algorithm 3.2 proceeds increasing the value of the regularization parameter $\mu \geq 0$ until the sufficient descent condition (8) is satisfied with $s = s(\mu)$. For each value of μ , we define $\rho(\mu) = ([-\lambda_{1,k}]_+ + \mu)/(3\|s(\mu)\|)$, where $s(\mu)$ is a solution to (9). By Lemma 3.1 of [5] (see also [15, 20]), $s(\mu)$ is a global minimizer of $\langle g^k, s \rangle + \frac{1}{2} s^T H_k s + \rho(\mu) \|s\|^3$. The way in which μ is increased is determined by two necessities related to $\rho(\mu)$: (a) the initial $\rho(\mu)$ at each iteration should not be excessively small and (b) the final $\rho(\mu)$ should not be excessively big. Essentially, the technical manipulation of the quadratic regularization parameter μ in the algorithm is motivated but these two apparently conflicting objectives which are necessary to obtain the complexity results.

Algorithm 3.2.

Step 1. Let $\lambda_{1,k}$ be the smallest eigenvalue of H_k . Consider the linear system

$$[H_k + ([-\lambda_{1,k}]_+ + \mu)I]s = -g^k. \quad (10)$$

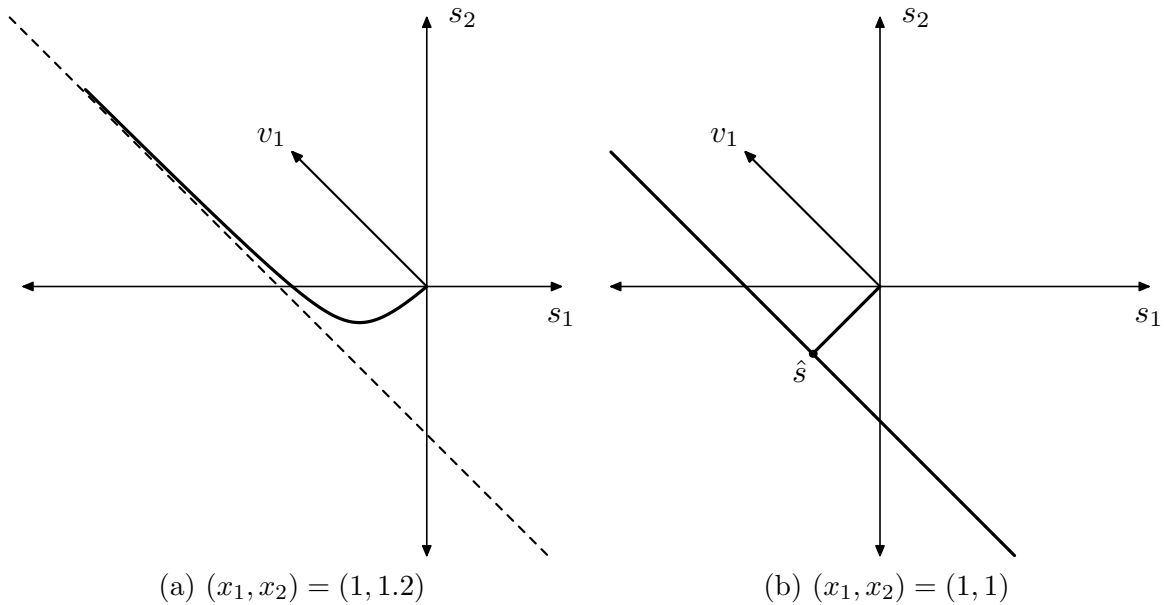


Figure 1: Levenberg-Marquardt path of the function $f(x_1, x_2) = x_1x_2$. The constant Hessian is indefinite, with eigenvalues $\lambda_1 = -1$ and $\lambda_2 = 1$. $v_1 = (-1, 1)$ is an eigenvector associated with λ_1 . The linear system (10) with $\mu = 0$ is given by $s_1 + s_2 = -x_2$ and $s_1 + s_2 = -x_1$, that is compatible if and only if $x_1 = x_2$. The picture in (a) corresponds to the case in which the linear system (10) with $\mu = 0$ is not compatible; while the picture in (b) represents the hard case and corresponds to the case in which the linear system is compatible.

If (10) with $\mu = 0$ is *not* compatible then set $\rho_{k,0} = 0$ and go to Step 5.

Step 2. Compute the minimum norm solution $\hat{s}^{k,0}$ to the linear system (10) with $\mu = 0$ and set

$$\rho_{k,0} = \begin{cases} \infty, & \text{if } \hat{s}^{k,0} = 0 \text{ and } [-\lambda_{1,k}]_+ > 0, \\ 0, & \text{if } \hat{s}^{k,0} = 0 \text{ and } [-\lambda_{1,k}]_+ = 0, \\ \frac{[-\lambda_{1,k}]_+}{3\|\hat{s}^{k,0}\|}, & \text{if } \hat{s}^{k,0} \neq 0. \end{cases}$$

If $\rho_{k,0} \leq M$ then go to Step 4.

Step 3. Let $q^{1,k}$ with $\|q^{1,k}\| = 1$ be an eigenvector of H_k associated with its smallest eigenvalue $\lambda_{1,k}$. Set $\ell_3 \leftarrow 1$ and compute $t_{\ell_3} \geq 0$ and $\hat{s}^{k,\ell_3} = \hat{s}^{k,0} + t_{\ell_3}q^{1,k}$ such that

$$\frac{[-\lambda_{1,k}]_+}{3\|\hat{s}^{k,\ell_3}\|} = M. \quad (11)$$

Step 3.1. If (8) holds with $s = \hat{s}^{k,\ell_3}$ then **return** $s = \hat{s}^{k,\ell_3}$.

Step 3.2. If

$$\|\hat{s}^{k,\ell_3}\| \geq 2\|\hat{s}^{k,0}\| \quad (12)$$

then compute $t_{\ell_3+1} \geq 0$ and $\hat{s}^{k,\ell_3+1} = \hat{s}^{k,0} + t_{\ell_3+1}q^{1,k}$ such that

$$\|\hat{s}^{k,\ell_3+1}\| = \frac{1}{2}\|\hat{s}^{k,\ell_3}\|, \quad (13)$$

set $\ell_3 \leftarrow \ell_3 + 1$, and go to Step 3.1.

Step 4. If (8) holds with $s = \hat{s}^{k,0}$ then **return** $s = \hat{s}^{k,0}$.

Step 5. Set $\ell_5 \leftarrow 1$ and $\rho_{k,\ell_5} = \max\{0.1, \rho_{k,0}\}$.

Step 5.1. Compute $\tilde{\mu}_{k,\ell_5} > 0$ and \tilde{s}^{k,ℓ_5} solution to (10) with $\mu = \tilde{\mu}_{k,\ell_5}$ such that

$$\rho_{k,\ell_5} \leq \frac{[-\lambda_{1,k}]_+ + \tilde{\mu}_{k,\ell_5}}{3\|\tilde{s}^{k,\ell_5}\|} \leq 100\rho_{k,\ell_5}. \quad (14)$$

Step 5.2. If (8) holds with $s = \tilde{s}^{k,\ell_5}$ then **return** $s = \tilde{s}^{k,\ell_5}$.

Step 5.3. If $\tilde{\mu}_{k,\ell_5} < 0.1$ then set

$$\rho_{k,\ell_5+1} = 10 \left(\frac{[-\lambda_{1,k}]_+ + \tilde{\mu}_{k,\ell_5}}{3\|\tilde{s}^{k,\ell_5}\|} \right) \quad (15)$$

and $\ell_5 \leftarrow \ell_5 + 1$, and go to Step 5.1.

Step 6. Set $\ell_6 \leftarrow 1$ and $\bar{\mu}_{k,\ell_6} = 2\tilde{\mu}_{k,\ell_5}$.

Step 6.1. Compute \bar{s}^{k,ℓ_6} solution to (10) with $\mu = \bar{\mu}_{k,\ell_6}$.

Step 6.2. If (8) holds with $s = \bar{s}^{k,\ell_6}$ then **return** $s = \bar{s}^{k,\ell_6}$.

Step 6.3. Set $\bar{\mu}_{k,\ell_6+1} = 2\bar{\mu}_{k,\ell_6}$ and $\ell_6 \leftarrow \ell_6 + 1$ and go to Step 6.1.

Let us comment how Algorithm 3.2 proceeds in each of the three possible cases (i) H_k is positive definite or H_k is positive semidefinite and the linear system (10) with $\mu = 0$ is compatible, (ii) H_k is *not* positive semidefinite and the linear system (10) with $\mu = 0$ is compatible, and (iii) H_k is *not* positive semidefinite and the linear system (10) with $\mu = 0$ is *not* compatible, that correspond to the three possible geometries of the Levenberg-Marquardt path.

If H_k is positive definite, we have that $\lambda_{1,k} > 0$, so $[-\lambda_{1,k}]_+ = 0$. Then, the system (10) with $\mu = 0$ is compatible and, by Step 2, $\rho_{k,0} = 0$. Since $\rho_{k,0} \leq M$, the algorithm continues at Step 4 and the increment $\hat{s}^{k,0}$ is accepted if the sufficient descent condition (8) holds with $s = \hat{s}^{k,0}$ (this is always the case if $\hat{s}^{k,0} = 0$ that occurs if and only if $g^k = 0$). However, if (8) does not hold, the algorithm computes at Step 5.1 a regularization parameter μ such that the corresponding $\rho(\mu)$

increases with respect to the previous one, but not very much. This corresponds to our purpose of maintaining the auxiliary quantity $\rho(\mu)$ within controlled bounds. If $s(\mu)$ does not satisfy (8) (Step 5.2) and the regularization parameter μ is still small (Step 5.3), we update (increase) the bounds on $\rho(\mu)$, and we repeat this process until the fulfillment of (8) or until μ is not small anymore. From this point on, the process continues with regular increases of the regularization parameter μ which should lead to the final fulfillment of (8) at Step 6.2. It is easy to see that, when H_k is positive semidefinite and the linear system $H_k s = -g^k$ is compatible, the algorithm proceeds as in the positive definite case described above.

If H_k is not positive semidefinite we have that $\lambda_{1,k} < 0$ and $[-\lambda_{1,k}]_+ > 0$. If, in spite of that, the system (10) with $\mu = 0$ is compatible, we are in presence of the hard case. In this case, we compute the minimum norm solution of (10), which corresponds to the intersection of the two branches of the Levenberg-Marquardt path in Figure 1b. If taking the regularizing parameter $\mu = 0$ we have that the associated $\rho(\mu)$ is not very big ($\rho_{k,0} \leq M$ at Step 2) then we proceed exactly as in the positive definite and compatible positive semidefinite cases, increasing μ and seeking an acceptable increment along the “vertical” branch of the Levenberg-Marquardt path (that joins the minimum-norm solution to the origin). However, if $\rho_{k,0} > M$, we are in the case that $\rho(\mu)$ could be very big. Then, it is time to seek an increment along the “horizontal” branch of the Levenberg-Marquardt path (line generated by adding a multiple of an eigenvector associated with $\lambda_{1,k}$ to the minimum-norm solution). This is the case when $\lambda_{1,k} < 0$ and $\hat{s}^{k,0} = 0$ (because $g^k = 0$), since in that case, we set $\rho_{k,0} = \infty$ at Step 2. Note that, along this branch, the value of $\mu = 0$ does not change and the reduction of $\rho(\mu)$ is achieved trivially increasing the norm of $s(\mu)$. Starting with a sufficiently large $\|s(\mu)\|$, and by means of successive reductions of $\|s(\mu)\|$, we seek the fulfillment of (8). However, after a finite number of reductions of $\|s(\mu)\|$ this norm becomes smaller than the minimum norm in the horizontal branch (except in the case in which we have $\hat{s}^{k,0} = 0$). If this happens, we initiate a search in the vertical branch in an analogous way as we do in the positive definite case. In this case, we have the guarantee that $\rho(\mu)$ is suitable bounded in the intersection point because, otherwise, the sufficient descent condition (8) would have been accepted.

Finally, let us consider the case in which H_k is not positive definite and the system (10) with $\mu = 0$ is not compatible. This situation corresponds to Figure 1a. In this case, the control goes to Step 5. At Steps 5.1, 5.2, and 5.3, the regularization parameter μ is increased in a controlled way, so that abrupt changes of the associated $\rho(\mu)$ are not possible. In this process, when μ becomes greater than or equal to 0.1, the increase-regularization process continues at Step 6 with a single doubling procedure.

4 Well-definiteness results

In this section, it will be shown that Algorithm 3.2 is well-defined and that the computed increment s^k that satisfies (8) also satisfies (3). We start by describing how Algorithm 3.2 could be implemented considering the spectral decomposition of H_k . Of course, this is an arbitrary choice and other options are possible. In any case, this description introduces some useful notation for the remaining of the section.

Consider the spectral decomposition

$$H_k = Q_k \Lambda_k Q_k^T,$$

where $\Lambda_k = \text{diag}(\lambda_{1,k}, \dots, \lambda_{n,k})$ with $\lambda_{1,k} \leq \dots \leq \lambda_{n,k}$ and $Q = [q^{1,k} \dots q^{n,k}]$ is orthogonal. Substituting H_k by its spectral decomposition in (10), we obtain

$$[\Lambda_k + ([-\lambda_{1,k}]_+ + \mu)I]Q_k^T s = -Q_k^T g^k.$$

Therefore, for $\mu = 0$, the linear system (10) is compatible if and only if $[Q_k^T g^k]_j = 0$ whenever

$$\lambda_{j,k} + [-\lambda_{1,k}]_+ = 0. \quad (16)$$

Assuming that the linear system (10) with $\mu = 0$ is compatible, its minimum norm solution is given by $\hat{s}^{k,0} = Q_k y^k$, where

$$y_j^k = \begin{cases} -[Q_k^T g^k]_j / (\lambda_{j,k} + [-\lambda_{1,k}]_+), & j \in J, \\ 0, & j \in \bar{J}, \end{cases}$$

$J = \{j \in \{1, \dots, n\} \mid \lambda_{j,k} + [-\lambda_{1,k}]_+ \neq 0\}$, and $\bar{J} = \{1, \dots, n\} \setminus J$. Moreover, note that

$$\|\hat{s}^{k,0}\| = \sqrt{\sum_{j \in J} ([Q_k^T g^k]_j / (\lambda_{j,k} + [-\lambda_{1,k}]_+))^2}.$$

The norm of $\hat{s}^{k,\ell_3} = \hat{s}^{k,0} + t_{\ell_3} q^{1,k}$ (for any $\ell_3 \geq 1$) computed at Step 3 is given by

$$\|\hat{s}^{k,\ell_3}\| = \sqrt{\|\hat{s}^{k,0}\|^2 + t_{\ell_3}^2 \langle \hat{s}^{k,0}, q^{1,k} \rangle + t_{\ell_3}^2}.$$

Thus, given a desired norm c_{ℓ_3} for \hat{s}^{k,ℓ_3} ($c_{\ell_3} = [-\lambda_{1,k}]_+ / (3M)$ when $\ell_3 = 1$ and $c_{\ell_3} = \frac{1}{2} \|\hat{s}^{k,\ell_3-1}\|$ when $\ell_3 > 1$), we have that

$$t_{\ell_3} = -\langle \hat{s}^{k,0}, q^{1,k} \rangle + \sqrt{\langle \hat{s}^{k,0}, q^{1,k} \rangle^2 - (\|\hat{s}^{k,0}\|^2 - c_{\ell_3})}.$$

The following technical lemma establishes that Step 5.1 of Algorithm 3.2 can always be completed finding a regularization parameter μ and an increment $s(\mu)$ that satisfies (14). The assumption $g^k \neq 0$ in the lemma is perfectly reasonable because, as it will be shown later, it always holds at Step 5.1.

Lemma 4.1 *Suppose that $g^k \neq 0$. At Step 5.1 of Algorithm 3.2, for any $\ell_5 \geq 1$, there exists $\tilde{\mu}_{k,\ell_5} > 0$ and \tilde{s}^{k,ℓ_5} solution to (10) with $\mu = \tilde{\mu}_{k,\ell_5}$ satisfying (14).*

Proof: For any $\mu > 0$, the matrix of the system (10) is positive definite and the solution $s(\mu)$ to (10) is such that

$$\|s(\mu)\| = \sqrt{\sum_{\{j \mid [Q_k^T g^k]_j \neq 0\}} \left(\frac{[Q_k^T g^k]_j}{(\lambda_{j,k} + [-\lambda_{1,k}]_+ + \mu)} \right)^2}. \quad (17)$$

Moreover, clearly,

$$\lim_{\mu \rightarrow \infty} \|s(\mu)\| = 0. \quad (18)$$

In order to analyze the case $\mu \rightarrow 0$, the proof will be divided in two cases: (a) the linear system (10) with $\mu = 0$ is compatible and (b) the linear system (10) with $\mu = 0$ is *not* compatible.

Consider first case (a). In this case, since $[Q_k^T g^k]_j = 0$ whenever $\lambda_{j,k} + [-\lambda_{1,k}]_+ = 0$, (17) is equivalent to

$$\|s(\mu)\| = \sqrt{\sum_{j \in J} \left(\frac{[Q_k^T g^k]_j}{(\lambda_{j,k} + [-\lambda_{1,k}]_+ + \mu)} \right)^2}.$$

Therefore,

$$\lim_{\mu \rightarrow 0} \|s(\mu)\| = \|\hat{s}^{k,0}\| > 0 \quad (19)$$

because $g^k \neq 0$ implies $\hat{s}^{k,0} \neq 0$. Thus, by (18) and (19), we have that

$$\lim_{\mu \rightarrow \infty} \frac{[-\lambda_{1,k}]_+ + \mu}{3\|s(\mu)\|} = \infty \quad \text{and} \quad \lim_{\mu \rightarrow 0} \frac{[-\lambda_{1,k}]_+ + \mu}{3\|s(\mu)\|} = \frac{[-\lambda_{1,k}]_+}{3\|\hat{s}^{k,0}\|}. \quad (20)$$

Since, by definition, for any $\ell_5 \geq 1$,

$$\rho_{k,\ell_5} \geq \rho_{k,0} = \frac{[-\lambda_{1,k}]_+}{3\|\hat{s}^{k,0}\|},$$

the desired result follows by continuity from (20).

Consider now case (b). In this case, there exists j such that $\lambda_{j,k} + [-\lambda_{1,k}]_+ = 0$ and $[Q_k^T g^k]_j \neq 0$. Therefore, from (17), we have that

$$\lim_{\mu \rightarrow 0} \|s(\mu)\| = \infty. \quad (21)$$

Thus, by (18) and (21), we have that

$$\lim_{\mu \rightarrow \infty} \frac{[-\lambda_{1,k}]_+ + \mu}{3\|s(\mu)\|} = \infty \quad \text{and} \quad \lim_{\mu \rightarrow 0} \frac{[-\lambda_{1,k}]_+ + \mu}{3\|s(\mu)\|} = 0. \quad (22)$$

Since, by definition, for any $\ell_5 \geq 1$, in this case we have

$$\rho_{k,\ell_5} \geq \rho_{k,0} = 0.1,$$

the desired result follows by continuity from (22). \square

Below we state the main assumption that supports the complexity results. Essentially, we will assume that the objective function is twice continuously differentiable and that $\nabla^2 f$ satisfies a Lipschitz condition on a suitable region that contains the iterates x^k and the trial points $x^k + s^{\text{trial}}$. Of course, a sufficient condition for the fulfillment of this assumption is the Lipschitz-continuity of $\nabla^2 f$ onto \mathbb{R}^n , but in some cases this global assumption may be unnecessarily strong.

Assumption A1 *The function f is twice continuous differentiable for all $x \in \mathbb{R}^n$ and there exists a constant $L > 0$ such that, for all x^k computed by Algorithm 3.1 and every trial increment s^{trial} computed at Steps 2, 3, 3.2, 5.1, or 6.1 of Algorithm 3.2, we have that*

$$f(x^k + s^{\text{trial}}) \leq f(x^k) + (s^{\text{trial}})^T g^k + \frac{1}{2}(s^{\text{trial}})^T H_k s^{\text{trial}} + L \|s^{\text{trial}}\|^3$$

and

$$\|g(x^k + s^{\text{trial}}) - g^k - H_k s^{\text{trial}}\| \leq L \|s^{\text{trial}}\|^2.$$

In the following lemma we prove that any trial increment necessarily satisfies the sufficient descent condition (8) if the regularization parameter is large enough.

Lemma 4.2 *Suppose that Assumption A1 holds and $\mu \geq 0$. If $0 \neq s^{\text{trial}} \in \mathbb{R}^n$ computed at Steps 2, 3, 3.2, 5.1, or 6.1 of Algorithm 3.2, that by definition satisfies*

$$[H_k + ([-\lambda_{1,k}]_+ + \mu)] s^{\text{trial}} = -g^k, \quad (23)$$

is such that

$$\frac{[-\lambda_{1,k}]_+ + \mu}{3 \|s^{\text{trial}}\|} \geq L + \alpha \quad (24)$$

then (8) is satisfied with $s = s^{\text{trial}}$.

Proof: Let us define, for all $s \in \mathbb{R}^n$,

$$q(s) = s^T g^k + \frac{1}{2} s^T H_k s.$$

Since $H_k + ([-\lambda_{1,k}]_+ + \mu)I$ is positive semidefinite for any $\mu \geq 0$, by (23),

$$s^{\text{trial}} \text{ minimizes } q(s) + \frac{1}{2} ([-\lambda_{1,k}]_+ + \mu) \|s\|^2. \quad (25)$$

Define

$$\rho = \frac{[-\lambda_{1,k}]_+ + \mu}{3 \|s^{\text{trial}}\|}. \quad (26)$$

By Lemma 3.1 of [5], s^{trial} is a minimizer of $q(s) + \rho \|s\|^3$. In particular,

$$q(s^{\text{trial}}) + \rho \|s^{\text{trial}}\|^3 \leq q(0) = 0. \quad (27)$$

Now, by Assumption A1, we have that

$$\begin{aligned} f(x^k + s^{\text{trial}}) &\leq f(x^k) + (s^{\text{trial}})^T g^k + \frac{1}{2} (s^{\text{trial}})^T H_k s^{\text{trial}} + L \|s^{\text{trial}}\|^3 \\ &= f(x^k) + q(s^{\text{trial}}) + \rho \|s^{\text{trial}}\|^3 + (L - \rho) \|s^{\text{trial}}\|^3. \end{aligned}$$

Thus, by (24), (26), and (27),

$$f(x^k + s^{\text{trial}}) \leq f(x^k) - \alpha \|s^{\text{trial}}\|^3.$$

This completes the proof. \square

The lemma below shows that Algorithm 3.2 may return a null increment only at Step 4.

Lemma 4.3 *Suppose that Assumption A1 holds. Algorithm 3.2 returns a null increment $s = 0$ if and only if $g^k = 0$ and $\lambda_{1,k} \geq 0$. Moreover, an increment $s = 0$ may only be returned by Algorithm 3.2 at Step 4 (i.e. Steps 3.1, 5.2, and 6.2 always return non null increments).*

Proof: Assume that $g^k = 0$ and $\lambda_{1,k} \geq 0$. Then, we have that the minimum norm solution $\hat{s}^{k,0}$ to the linear system (10) with $\mu = 0$ computed at Step 2 is null and that $\rho_{k,0} = 0 \leq M$. Therefore, the algorithm goes to Step 4 and returns $s = \hat{s}^{k,0} = 0$ since it satisfies (8).

Assume now that Algorithm 3.2 returned an increment $s = 0$. Since every trial increment computed by the algorithm is a solution to the linear system (10) for some $\mu \geq 0$, we must have $g^k = 0$. If $\lambda_{1,k} \geq 0$ then the first part of thesis holds and it remains to show that the null increment is returned at Step 4. Note that, since $g^k = 0$ implies $\hat{s}^{k,0} = 0$ and $\lambda_{1,k} \geq 0$ means $[-\lambda_{1,k}]_+ = 0$, at Step 2 we have $\rho_{k,0} = 0 \leq M$. Thus, the algorithm goes to Step 4 where the null increment is returned since it satisfies (8). We now show that assuming $\lambda_{1,k} < 0$ leaves to a contradiction. Since $\lambda_{1,k} < 0$ means $[-\lambda_{1,k}]_+ > 0$ and $g^k = 0$ implies $\hat{s}^{k,0} = 0$, by the way $\rho_{k,0}$ is defined at Step 2, we have that $\rho_{k,0} = \infty \not\leq M$. In this case the algorithm goes to Step 3. On the one hand, note that $\hat{s}^{k,0} = 0$ implies that the algorithm never leaves the loop given by Steps 3.1–3.2 because (12) reduces to $\|\hat{s}^{k,\ell_3}\| \geq 0$. On the other hand, note that, by halving the norm of the trial increments \hat{s}^{k,ℓ_3} , since $\mu = 0$ is fixed, in a finite number of trials, (24) holds and, by Lemma 4.2, the algorithm returns $s = \hat{s}^{k,\ell_3} \neq 0$ for some $\ell_3 \geq 1$, contradicting the fact that the algorithm returned a null increment. \square

We finish this section proving that the increment s^k computed at Algorithm 3.2, that satisfies (8) and defines x^{k+1} in Algorithm 3.1, is such that it also satisfies (3). Note that this result assumes the existence of s^k by hypothesis. Up to the present moment we proved that Algorithm 3.2 is well defined. The existence of s^k for all k will be proved in the following section when proving that Algorithm 3.2 always computes s^k performing a finite number of operations.

Lemma 4.4 *Suppose that Assumption A1 holds. Then, there exist $\gamma_g > 0$ and $\gamma_H > 0$ such that, for all $k \in \mathbb{N}$, the increment s^k computed by Algorithm 3.2 and the new iterate $x^{k+1} = x^k + s^k$ computed at Step 2 of Algorithm 3.1 satisfy*

$$\sqrt{\frac{\|g^{k+1}\|}{\gamma_g}} \leq \|s^k\| \quad \text{and} \quad \frac{[-\lambda_{1,k}]_+}{\gamma_H} \leq \|s^k\|.$$

Moreover,

$$\gamma_g \leq \max \{3M + L, 3000(L + \alpha) + L, 30 + L\} \quad (28)$$

and

$$\gamma_H \leq \max \{3M, 3000(L + \alpha), 30\}. \quad (29)$$

Proof: If $s^k = 0$ then, by Lemma 4.3, we have that $g^k = 0$ and $\lambda_{1,k} \geq 0$ and, therefore, the thesis follows trivially. We now assume $s^k \neq 0$. Since s^k is a solution to (10) for some $\mu \geq 0$, we have that

$$H_k s^k + g^k + ([-\lambda_{1,k}]_+ + \mu) s^k = 0.$$

Therefore,

$$H_k s^k + g^k + \left(\frac{[-\lambda_{1,k}]_+ + \mu}{\|s^k\|} \right) \|s^k\| s^k = 0.$$

Then

$$\|H_k s^k + g^k\| = \left(\frac{[-\lambda_{1,k}]_+ + \mu}{\|s^k\|} \right) \|s^k\|^2.$$

But, by Assumption A1 and the triangle inequality,

$$\|g^{k+1}\| - \|g^k + H_k s^k\| \leq \|g^{k+1} - g^k - H_k s^k\| \leq L \|s^k\|^2.$$

Therefore,

$$\|g^{k+1}\| \leq \left(\frac{[-\lambda_{1,k}]_+ + \mu}{\|s^k\|} + L \right) \|s^k\|^2. \quad (30)$$

We now analyze in separate the cases in which $s^k \neq 0$ is returned by Algorithm 3.2 at Steps 3.1, 4, 5.2, and 6.2.

Case $s^k = \hat{s}^{k,\ell_3}$ with $\ell_3 = 1$ was returned at Step 3.1: In this case, s^{k,ℓ_3} is a solution to (10) with $\mu = 0$ and, by (11), it satisfies

$$\frac{[-\lambda_{1,k}]_+}{\|s^{k,\ell_3}\|} = 3M. \quad (31)$$

Case $s^k = \hat{s}^{k,\ell_3}$ with $\ell_3 > 1$ was returned at Step 3.1: This means that there exists $\hat{s}^{k,\ell_3-1} \neq 0$ that is a solution to (10) with $\mu = 0$ and for which (8) with $s = \hat{s}^{k,\ell_3-1}$ did not hold. Therefore, by Lemma 4.2, we have that

$$\frac{[-\lambda_{1,k}]_+}{3\|\hat{s}^{k,\ell_3-1}\|} < L + \alpha.$$

Thus, by (13), we have that

$$\frac{[-\lambda_{1,k}]_+}{\|\hat{s}^{k,\ell_3}\|} < 6(L + \alpha). \quad (32)$$

Case $s^k = \hat{s}^{k,0}$ was returned at Step 4: In this case, we have that

$$\frac{[-\lambda_{1,k}]_+}{3\|\hat{s}^{k,0}\|} \leq M \quad (33)$$

or that there exists $\hat{s}^{k,\ell_3} \neq 0$ with $\ell_3 \geq 1$ such that

$$\|\hat{s}^{k,\ell_3}\| < 2\|\hat{s}^{k,0}\|, \quad (34)$$

\hat{s}^{k,ℓ_3} is a solution to (10) with $\mu = 0$, and (8) did not hold with $s = \hat{s}^{k,\ell_3}$. Therefore, by Lemma 4.2, we have that

$$\frac{[-\lambda_{1,k}]_+}{3\|\hat{s}^{k,\ell_3}\|} < L + \alpha \quad (35)$$

and, by (34,35),

$$\frac{[-\lambda_{1,k}]_+}{\|\tilde{s}^{k,0}\|} < 6(L + \alpha). \quad (36)$$

Thus, by (33) and (36),

$$\frac{[-\lambda_{1,k}]_+}{\|s^{k,0}\|} \leq \max\{3M, 6(L + \alpha)\}. \quad (37)$$

Case $s^k = \tilde{s}^{k,\ell_5}$ with $\ell_5 = 1$ was returned at Step 5.2: In this case there are two possibilities: the linear system (10) with $\mu = 0$ is compatible or not. In the first case, $\hat{s}^{k,0}$ was computed,

$$\rho_{k,0} = \frac{[-\lambda_{1,k}]_+}{3\|\hat{s}^{k,0}\|},$$

and, since (8) with $s = \hat{s}^{k,0}$ did not hold, by Lemma 4.2, $\rho_{k,0} < L + \alpha$. In the second case, we simple have that $\rho_{k,0} = 0$. Thus, by (14) and by the fact that, by definition,

$$\rho_{k,1} = \max\{0.1, \rho_{k,0}\},$$

in the first case, we have

$$\frac{[-\lambda_{1,k}]_+ + \tilde{\mu}_{k,\ell_5}}{3\|\tilde{s}^{k,\ell_5}\|} \leq 100\rho_{k,\ell_5} = 100 \max\{0.1, \rho_{k,0}\} \leq \max\{10, 100(L + \alpha)\} \quad (38)$$

and, in the second case, we have

$$\frac{[-\lambda_{1,k}]_+ + \tilde{\mu}_{k,\ell_5}}{3\|\tilde{s}^{k,\ell_5}\|} \leq 100\rho_{k,\ell_5} = 100 \max\{0.1, 0\} = 10. \quad (39)$$

Therefore, $\tilde{\mu}_{k,\ell_5} \geq 0$, (38), and (39) imply that

$$\frac{[-\lambda_{1,k}]_+}{\|s^{k,\ell_5}\|} \leq \frac{[-\lambda_{1,k}]_+ + \tilde{\mu}_{k,\ell_5}}{\|s^{k,\ell_5}\|} \leq \max\{30, 300(L + \alpha)\}. \quad (40)$$

Case $s^k = \tilde{s}^{k,\ell_5}$ with $\ell_5 > 1$ was returned at Step 5.2: This means that there exists $\tilde{\mu}_{k,\ell_5-1} > 0$ and \tilde{s}^{k,ℓ_5-1} solution to (10) with $\mu = \tilde{\mu}_{k,\ell_5-1}$ for which (8) did not hold. Thus, by Lemma 4.2,

$$\frac{[-\lambda_{1,k}]_+ + \tilde{\mu}_{k,\ell_5-1}}{3\|\tilde{s}^{k,\ell_5-1}\|} < L + \alpha.$$

Moreover, by (14) and (15),

$$\frac{[-\lambda_{1,k}]_+ + \tilde{\mu}_{k,\ell_5}}{3\|\tilde{s}^{k,\ell_5}\|} \leq 100\rho_{k,\ell_5} = 1000 \left(\frac{[-\lambda_{1,k}]_+ + \tilde{\mu}_{k,\ell_5-1}}{3\|\tilde{s}^{k,\ell_5-1}\|} \right).$$

Thus,

$$\frac{[-\lambda_{1,k}]_+}{\|\tilde{s}^{k,\ell_5}\|} \leq \frac{[-\lambda_{1,k}]_+ + \tilde{\mu}_{k,\ell_5}}{\|\tilde{s}^{k,\ell_5}\|} \leq 3000(L + \alpha). \quad (41)$$

Case $s^k = \bar{s}^{k,\ell_6}$ was returned at Step 6.2: If $\ell_6 = 1$ then $\bar{\mu}_{k,\ell_6} = 2\tilde{\mu}_{k,\ell_5}$ for some $\ell_5 \geq 1$ and the solution \tilde{s}^{k,ℓ_5} to (10) with $\mu = \tilde{\mu}_{k,\ell_5}$ is such that (8) with $s = \tilde{s}^{k,\ell_5}$ does not hold. Thus, by Lemma 4.2,

$$\frac{[-\lambda_{1,k}]_+ + \tilde{\mu}_{k,\ell_5}}{3\|\tilde{s}^{k,\ell_5}\|} < L + \alpha.$$

On the other, since $\bar{\mu}_{k,\ell_6} = 2\tilde{\mu}_{k,\ell_5}$, we have that

$$\begin{aligned} \|\bar{s}^{k,\ell_6}\| &= \sqrt{\sum_{j \in J} \left(\frac{[Q_k^T g^k]_j}{\lambda_{j,k} + [-\lambda_{1,k}]_+ + \bar{\mu}_{k,\ell_6}} \right)^2} = \sqrt{\sum_{j \in J} \left(\frac{[Q_k^T g^k]_j}{\lambda_{j,k} + [-\lambda_{1,k}]_+ + 2\tilde{\mu}_{k,\ell_5}} \right)^2} \\ &= \sqrt{\sum_{j \in J} \left(\frac{[Q_k^T g^k]_j}{2(\frac{1}{2}(\lambda_{j,k} + [-\lambda_{1,k}]_+) + \tilde{\mu}_{k,\ell_5})} \right)^2} \geq \sqrt{\sum_{j \in J} \left(\frac{[Q_k^T g^k]_j}{2(\lambda_{j,k} + [-\lambda_{1,k}]_+ + \tilde{\mu}_{k,\ell_5})} \right)^2} \\ &= \frac{1}{2} \sqrt{\sum_{j \in J} \left(\frac{[Q_k^T g^k]_j}{\lambda_{j,k} + [-\lambda_{1,k}]_+ + \tilde{\mu}_{k,\ell_5}} \right)^2} = \frac{1}{2} \|\tilde{s}^{k,\ell_5}\| > 0. \end{aligned} \tag{42}$$

Therefore,

$$\begin{aligned} \frac{[-\lambda_{1,k}]_+}{\|\bar{s}^{k,\ell_6}\|} &\leq \frac{[-\lambda_{1,k}]_+ + \bar{\mu}_{k,\ell_6}}{\|\bar{s}^{k,\ell_6}\|} = \frac{[-\lambda_{1,k}]_+ + 2\tilde{\mu}_{k,\ell_5}}{\|\bar{s}^{k,\ell_6}\|} = \frac{2(\frac{1}{2}[-\lambda_{1,k}]_+ + \tilde{\mu}_{k,\ell_5})}{\|\bar{s}^{k,\ell_6}\|} \leq \\ &\frac{2([- \lambda_{1,k}]_+ + \tilde{\mu}_{k,\ell_5})}{\|\bar{s}^{k,\ell_6}\|} \leq \frac{2([- \lambda_{1,k}]_+ + \tilde{\mu}_{k,\ell_5})}{\frac{1}{2}\|\tilde{s}^{k,\ell_5}\|} = 4 \left(\frac{[-\lambda_{1,k}]_+ + \tilde{\mu}_{k,\ell_5}}{\|\tilde{s}^{k,\ell_5}\|} \right) < 12(L + \alpha). \end{aligned} \tag{43}$$

If $\ell_6 > 1$ then $\bar{\mu}_{k,\ell_6} = 2\bar{\mu}_{k,\ell_6-1}$ and the solution \bar{s}^{k,ℓ_6-1} to (10) with $\mu = \bar{\mu}_{k,\ell_6-1}$ is such that (8) with $s = \bar{s}^{k,\ell_6-1}$ does not hold. Thus, by Lemma 4.2,

$$\frac{[-\lambda_{1,k}]_+ + \bar{\mu}_{k,\ell_6-1}}{3\|\bar{s}^{k,\ell_6-1}\|} < L + \alpha. \tag{44}$$

Moreover, $\bar{\mu}_{k,\ell_6} = 2\bar{\mu}_{k,\ell_6-1}$ implies, as shown above, that

$$\|\bar{s}^{k,\ell_6}\| \geq \frac{1}{2} \|\bar{s}^{k,\ell_6-1}\|. \tag{45}$$

Therefore, by (44) and (45), and since $\bar{\mu}_{k,\ell_6} \geq 0$, we have that

$$\frac{[-\lambda_{1,k}]_+}{\|\bar{s}^{k,\ell_6}\|} \leq \frac{[-\lambda_{1,k}]_+ + \bar{\mu}_{k,\ell_6}}{\|\bar{s}^{k,\ell_6}\|} < 12(L + \alpha). \tag{46}$$

The desired result (29) follows from (31,32,37,40,41,43,46); while (28) follows from the same set of inequalities plus (30). \square

5 Complexity results

In this section, complexity results on Algorithm 3.2 are presented. In particular, we show that the number of functional evaluations required to compute the increment s^k using Algorithm 3.2 is $O(1)$, i.e. it does not depend on ε_g nor ε_H . The section finishes establishing the complexity of Algorithm 3.1–3.2 in terms of the number of functional (and derivatives) evaluations. The sufficient condition (8) is tested at Steps 3.1, 4, 5.2, and 6.2. These are the only steps of Algorithm 3.2 in which the objective function is evaluated. Therefore, in order to assess the worst-case evaluation complexity of Algorithm 3.2, we must obtain a bound for the number of executions of each of these steps.

Step 3.1 corresponds to the hard case, in which we seek an increment along an appropriate eigenvector of H_k . Along this process, either an increment is accepted or the process finishes with Step 4, where f is evaluated (at most once per iteration) at the bifurcation point of Figure 1b. The lemma below establishes an upper bound on the number of executions of Step 3.1.

Lemma 5.1 *Suppose that Assumption A1 holds. If Step 3.1 of Algorithm 3.2 is executed, it is executed at most $\lfloor \log_2((L + \alpha)/M) \rfloor + 1$ times.*

Proof: By construction, $\hat{s}^{k,\ell_3} \neq 0$ for all $\ell_3 \geq 0$ and

$$\|\hat{s}^{k,\ell_3}\| = \begin{cases} [-\lambda_{1,k}]_+ / (3M), & \ell_3 = 1, \\ \|\hat{s}^{k,\ell_3-1}\| / 2, & \ell_3 > 1, \end{cases}$$

or, equivalently,

$$2^{\ell_3-1} M = \frac{[-\lambda_{1,k}]_+}{3\|\hat{s}^{k,\ell_3}\|}. \quad (47)$$

Thus, by Lemma 4.2, if (8) does not hold with $s = \hat{s}^{k,\ell_3}$ we must have $2^{\ell_3-1} M < L + \alpha$, i.e. $\ell_3 \leq \lfloor \log_2((L + \alpha)/M) \rfloor + 1$ as we wanted to prove. \square

Step 5 of Algorithm 3.2 describes a loop where one tries to find an “initial” sufficiently big regularization parameter. Each time the regularization parameter is increased one tests the condition (8). Therefore, it is necessary to establish a bound on the number of evaluations that may be performed at Step 5.2. This is done in Lemma 5.2.

Lemma 5.2 *Suppose that Assumption A1 holds. If Step 5.2 of Algorithm 3.2 is executed, it is executed at most $\lfloor \log_{10}(L + \alpha) \rfloor + 2$ times.*

Proof: For all $\ell_5 \geq 1$, when (8) is tested at Step 5.2 with $s = \tilde{s}^{k,\ell_5}$, \tilde{s}^{k,ℓ_5} is a solution to (10) with $\mu = \tilde{\mu}_{k,\ell_5} > 0$ and satisfies (14). Therefore, by Lemma 4.3, $\tilde{s}^{k,\ell_5} \neq 0$ and, thus, by Lemma 4.2, if (8) does not hold with $s = \tilde{s}^{k,\ell_5}$ we must have

$$\rho_{k,\ell_5} < L + \alpha. \quad (48)$$

On the other hand, since, by definition, $\rho_{k,1} \geq 0.1$ and, by (14) and (15), $\rho_{k,\ell_5} \geq 10\rho_{k,\ell_5-1}$ for all $\ell_5 \geq 2$, we have that

$$\rho_{k,\ell_5} \geq 10^{\ell_5-2} \quad (49)$$

for all $\ell_5 \geq 1$. By (48) and (49), if (8) does not hold with $s = \bar{s}^{k, \ell_5}$ we must have $10^{\ell_5 - 2} < L + \alpha$, i.e. $\ell_5 \leq \lfloor \log_{10}(L + \alpha) \rfloor + 2$ as we wanted to prove. \square

Finally, at Step 6 we increase the regularization parameter by means of a doubling process ($\bar{\mu}_{k, \ell_6 + 1} = 2\bar{\mu}_{k, \ell_6}$). This process guarantees, by Lemma 4.3 and Lemma 4.2, that the sufficient condition will eventually hold. In Lemma 5.3, we prove that the number of doubling steps is also bounded by a quantity that only depends on characteristics of the problem and algorithmic parameters. For proving this lemma, we need to assume boundedness of $\|H(x^k)\|$ at the iterates generated by the algorithm. Note that, since $f(x^{k+1}) \leq f(x^k)$ for all k , a sufficient condition for Assumption A2 is the boundedness of $\|H(x)\|$ on the level set defined by $f(x^0)$.

Assumption A2 *There exists a constant $h_{\max} \geq 0$ such that, for all iterate x^k computed by Algorithm 3.1, we have that $\|H(x^k)\| \leq h_{\max}$.*

Lemma 5.3 *Suppose that Assumptions A1 and A2 hold. If Step 6.2 of Algorithm 3.2 is executed, it is executed at most*

$$\left\lceil \left[\log \left(1 + \frac{0.2}{h_{\max} + 0.2} \right) \right]^{-1} \log \left(\frac{L + \alpha}{0.1} \right) \right\rceil + 1$$

times.

Proof: For all $\ell_6 \geq 1$, Lemma 4.3 implies that $\bar{s}^{k, \ell_6} \neq 0$ and straightforward calculations show that

$$\|\bar{s}^{k, \ell_6}\| = \sqrt{\sum_{j \in J} ([Q_k^T g^k]_j / (\lambda_{j, k} + [-\lambda_{1, k}]_+ + \bar{\mu}_{k, \ell_6}))^2}.$$

Moreover, it is easy to see that $\|\bar{s}^{k, \ell_6}\|$ decreases when $\bar{\mu}_{k, \ell_6}$ increases. Therefore, since, by definition, $\bar{\mu}_{k, \ell_6 + 1} = 2\bar{\mu}_{k, \ell_6}$, for all $\ell_6 \geq 1$, we have that

$$\frac{\|\bar{s}^{k, \ell_6}\|}{\|\bar{s}^{k, \ell_6 + 1}\|} \geq 1. \quad (50)$$

Thus, for all $\ell_6 \geq 1$,

$$\begin{aligned} & \left(\frac{[-\lambda_{1, k}]_+ + \bar{\mu}_{k, \ell_6 + 1}}{3\|\bar{s}^{k, \ell_6 + 1}\|} \right) / \left(\frac{[-\lambda_{1, k}]_+ + \bar{\mu}_{k, \ell_6}}{3\|\bar{s}^{k, \ell_6}\|} \right) = \left(\frac{[-\lambda_{1, k}]_+ + \bar{\mu}_{k, \ell_6 + 1}}{[-\lambda_{1, k}]_+ + \bar{\mu}_{k, \ell_6}} \right) \left(\frac{\|\bar{s}^{k, \ell_6}\|}{\|\bar{s}^{k, \ell_6 + 1}\|} \right) \geq \\ & \frac{[-\lambda_{1, k}]_+ + \bar{\mu}_{k, \ell_6 + 1}}{[-\lambda_{1, k}]_+ + \bar{\mu}_{k, \ell_6}} = \frac{[-\lambda_{1, k}]_+ + 2\bar{\mu}_{k, \ell_6}}{[-\lambda_{1, k}]_+ + \bar{\mu}_{k, \ell_6}} = 1 + \frac{\bar{\mu}_{k, \ell_6}}{[-\lambda_{1, k}]_+ + \bar{\mu}_{k, \ell_6}} \geq \left(1 + \frac{0.2}{h_{\max} + 0.2} \right) > 1, \end{aligned} \quad (51)$$

where the first inequality follows from (50) and the second inequality follows from the fact that, by the definition of the algorithm, $\bar{\mu}_{k, \ell_6} \geq 0.2$ and by Assumption A2.

From (51) and the fact that, by the definition of the algorithm, $\ell_6 = 1$ implies

$$\frac{[-\lambda_{1, k}]_+ + \bar{\mu}_{k, \ell_6}}{3\|\bar{s}^{k, \ell_6}\|} \geq 0.1,$$

it follows that

$$\frac{[-\lambda_{1,k}]_+ + \bar{\mu}_{k,\ell_6}}{3\|\bar{s}^{k,\ell_6}\|} \geq 0.1 \left(1 + \frac{0.2}{h_{\max} + 0.2}\right)^{\ell_6 - 1} \quad (52)$$

for all $\ell_6 \geq 1$. For all $\ell_6 \geq 1$, when (8) is tested at Step 6.2 with $s = \bar{s}^{k,\ell_6}$, \bar{s}^{k,ℓ_6} satisfies (10) with $\mu = \bar{\mu}_{k,\ell_6} > 0$. Therefore, by Lemma 4.2, if (8) does not hold with $s = \bar{s}^{k,\ell_6}$ we must have, by (52),

$$0.1 \left(1 + \frac{0.2}{h_{\max} + 0.2}\right)^{\ell_6 - 1} < L + \alpha.$$

This implies the desired result. \square

We finish this section summarizing the complexity and asymptotic results on Algorithm 3.1–3.2. Theorem 5.1 and Corollaries 5.1–5.3 presented below are analogous to Theorem 2.1 and Corollaries 2.1–2.3, respectively, presented for the model Algorithm 2.1 at the end of Section 2.

Theorem 5.1 *Let $f_{\min} \in \mathbb{R}$, $\varepsilon_g > 0$, and $\varepsilon_H > 0$ be given constants, suppose that Assumptions A1 and A2 hold, and let $\{x^k\}_{k=0}^\infty$ be the sequence generated by Algorithm 3.1–3.2. Then, the cardinality of the set of indices*

$$K_g = \left\{k \in \mathbb{N} \mid f(x^k) > f_{\min} \text{ and } \|g^{k+1}\| > \varepsilon_g\right\} \quad (53)$$

is, at most,

$$\left\lceil \frac{1}{\alpha} \left(\frac{f(x^0) - f_{\min}}{(\varepsilon_g/\gamma_g)^{3/2}} \right) \right\rceil; \quad (54)$$

while the cardinality of the set of indices

$$K_H = \left\{k \in \mathbb{N} \mid f(x^k) > f_{\min} \text{ and } \lambda_{1,k} < -\varepsilon_H\right\} \quad (55)$$

is, at most,

$$\left\lceil \frac{1}{\alpha} \left(\frac{f(x^0) - f_{\min}}{(\varepsilon_H/\gamma_H)^3} \right) \right\rceil, \quad (56)$$

where constants γ_g and γ_H are as in the thesis of Lemma 4.4 (i.e. they satisfy (28) and (29), respectively).

Proof: Assumptions A1 and A2 imply, by Lemma 4.4, the the hypothesis of Lemma 2.1 hold. Therefore, since Algorithm 3.1 is a particular case of Algorithm 2.1, from Lemma 2.1 it follows that at every time an iterate x^k is such that $\|g^{k+1}\| > \varepsilon_g$ the value of f decreases at least $\alpha(\varepsilon_g/\gamma_g)^{3/2}$; while at every time an iterate x^k is such that $\lambda_{1,k} < -\varepsilon_H$ the value of f decrease at least $\alpha(\varepsilon_H/\gamma_H)^3$. Therefore, the thesis follows from the fact that, by (8), $\{f(x^k)\}_{k=0}^\infty$ is a non-increasing sequence. \square

Corollary 5.1 *Let $f_{\min} \in \mathbb{R}$, $\varepsilon_g > 0$, and $\varepsilon_H > 0$ be given constants and suppose that Assumptions A1 and A2 hold. Algorithm 3.1–3.2 requires $O(\varepsilon_g^{-3/2})$ iterations and evaluations of f and its first- and second-order derivatives to compute x^k such that*

$$f(x^k) \leq f_{\min} \text{ or } \|g^{k+1}\| \leq \varepsilon_g;$$

it requires $O(\varepsilon_{\text{H}}^{-3})$ iterations and evaluations of f and its first- and second-order derivatives to compute x^k such that

$$f(x^k) \leq f_{\min} \text{ or } \lambda_{1,k} \geq -\varepsilon_{\text{H}};$$

and it requires $O(\varepsilon_g^{-3/2} + \varepsilon_{\text{H}}^{-3})$ iterations and evaluations of f and its first- and second-order derivatives to compute x^k such that

$$f(x^k) \leq f_{\min} \text{ or } \left(\|g^{k+1}\| \leq \varepsilon_g \text{ and } \lambda_{1,k} \geq -\varepsilon_{\text{H}} \right).$$

Proof: The required number of iterations is a direct consequence of Theorem 5.1. On the other hand, Lemmas 5.1–5.3 show that, every time Algorithm 3.2 is used by Algorithm 3.1 to compute an increment s^k , it performs $O(1)$ evaluations of the objective function f ; while, by definition, it performs a single evaluation of g and H . Thus, the evaluation complexity of Algorithm 3.1–3.2 coincides with its iteration complexity. \square

Corollary 5.2 *Suppose that Assumptions A1 and A2 hold and let $\{x^k\}_{k=0}^{\infty}$ be the sequence generated by Algorithm 3.1–3.2. Then, if the objective function f is bounded below, we have that*

$$\lim_{k \rightarrow \infty} \|g(x^k)\| = 0 \text{ and } \lim_{k \rightarrow \infty} [-\lambda_{1,k}]_+ = 0.$$

Proof: Assume that $\lim_{k \rightarrow \infty} \|g(x^k)\| \neq 0$. This means that there exists $\varepsilon > 0$ and \mathbb{K} , an infinite subsequence of \mathbb{N} , such that $\|g^k\| > \varepsilon$ for all $k \in \mathbb{K}$. Since f is bounded below, this contradicts Theorem 5.1. The second part is analogous. \square

Corollary 5.3 *Suppose that Assumptions A1 and A2 hold. Then, if the objective function f is bounded below, every limit point x^* of the sequence $\{x^k\}_{k=0}^{\infty}$ generated by Algorithm 3.1–3.2 is such that $\nabla f(x^*) = 0$ and $\nabla^2 f(x^*)$ is positive semidefinite.*

Proof: This corollary follows from Corollary 5.2 by continuity of ∇f and $\nabla^2 f$. \square

6 Local convergence

Note that if H_k is positive definite then the minimum norm solution $\hat{s}^{k,0}$ to the linear system (10) with $\mu = 0$ computed at Step 2 of Algorithm 3.2 is given by $\hat{s}^{k,0} = -H_k^{-1}g^k$, i.e. $\hat{s}^{k,0}$ is the Newton direction. Moreover, since, independently of having $\hat{s}^{k,0} = 0$ or $\hat{s}^{k,0} \neq 0$, $\lambda_{1,k} > 0$ implies that $\rho_{k,0} = 0 \leq M$, in this case (H_k positive definite) the algorithm goes directly to Step 4 and checks whether the Newton direction satisfies the sufficient cubic decrease condition (8). The lemma below shows that, if (57) holds then the Newton direction satisfies (8). (If $\lambda_{1,k} > 0$ and $g^k = 0$ and, in consequence, $s^{k,0} = 0$, it is trivial to see that the (null) Newton direction satisfies (8) and there is nothing to be proved. Anyway, the lemma below covers this case as well.)

Lemma 6.1 *Suppose that Assumption A1 holds. If H_k is positive definite and*

$$\|g^k\| \leq \frac{1}{2(L + \alpha)} \lambda_{1,k}^2 \tag{57}$$

then we have that the trial increment $\hat{s}^{k,0}$ computed at Step 2 of Algorithm 3.2 is such that (8) holds with $s = \hat{s}^{k,0}$.

Proof: By Assumption A1,

$$f(x^k + \hat{s}^{k,0}) \leq f(x^k) + (\hat{s}^{k,0})^T g^k + \frac{1}{2}(\hat{s}^{k,0})^T H_k \hat{s}^{k,0} + L\|\hat{s}^{k,0}\|^3.$$

Then, since $\hat{s}^{k,0} = -H_k^{-1}g^k$,

$$f(x^k + \hat{s}^{k,0}) \leq f(x^k) - \frac{1}{2}(\hat{s}^{k,0})^T H_k \hat{s}^{k,0} + L\|\hat{s}^{k,0}\|^3.$$

Therefore,

$$f(x^k + \hat{s}^{k,0}) \leq f(x^k) - \frac{1}{2}\lambda_{1,k}\|\hat{s}^{k,0}\|^2 + L\|\hat{s}^{k,0}\|^3. \quad (58)$$

On the other hand, since $\hat{s}^{k,0} = -H_k^{-1}g^k$, we have that

$$\|\hat{s}^{k,0}\| = \|H_k^{-1}g^k\| \leq \|H_k^{-1}\| \|g^k\| = \frac{1}{\lambda_{1,k}} \|g^k\|. \quad (59)$$

Then, by (57),

$$\|\hat{s}^{k,0}\| \leq \frac{\lambda_{1,k}}{2(L + \alpha)}$$

or, equivalently,

$$-\lambda_{1,k}/2 + L\|\hat{s}^{k,0}\| \leq -\alpha\|\hat{s}^{k,0}\|.$$

Therefore, multiplying by $\|\hat{s}^{k,0}\|^2$ and adding $f(x^k)$, we have that

$$f(x^k) - \frac{1}{2}\lambda_{1,k}\|\hat{s}^{k,0}\|^2 + L\|\hat{s}^{k,0}\|^3 \leq f(x^k) - \alpha\|\hat{s}^{k,0}\|^3.$$

and the thesis follows from (58). \square

In the next theorem, we use the classical local convergence result of Newton's method plus continuity arguments (that imply that the hypothesis (57) always hold in a neighborhood of a local minimizer with positive definite Hessian) to prove the quadratic local convergence of Algorithm 3.1–3.2.

Assumption A3 *Let x^* be a local minimizer of f . We say that this assumption holds if $H(x^*)$ is positive definite with $\|H(x^*)^{-1}\| \leq \beta$ and, in addition, there exist $r > 0$ and $\gamma > 0$ such that*

$$\|H(x) - H(x^*)\| \leq \gamma\|x - x^*\| \text{ whenever } \|x - x^*\| \leq r.$$

Theorem 6.1 *Let x^* be a local minimizer of f at which Assumption A3 holds and suppose that Assumption A1 also holds. Define $\delta_1 = \min\{r, \frac{1}{2\beta\gamma}\}$. Then, there exists $\delta \in (0, \delta_1]$ such that*

$$\|H(x)^{-1}\| \leq 2\beta \text{ whenever } \|x - x^*\| \leq \delta \quad (60)$$

and such that, if $\|x^0 - x^*\| \leq \delta$, the sequence $\{x^k\}_{k=0}^\infty$ generated by Algorithm 3.1–3.2 satisfies

$$\|g(x^k)\| \leq \left[\frac{1}{2(L + \alpha)} \right] / (2\beta)^2, \quad (61)$$

$$\|x^{k+1} - x^*\| \leq \frac{1}{2}\|x^k - x^*\|, \text{ and } \|x^{k+1} - x^*\| \leq \beta\gamma\|x^k - x^*\|^2 \quad (62)$$

for all $k = 0, 1, 2, \dots$.

Proof: By the classical Newton convergence theory (see, for example, [9, Th.5.2.1, p.90]), whenever $\|x^0 - x^*\| \leq \delta_1$ the sequence generated by $x^{k+1} = x^k - H_k^{-1}g^k$ is well defined and satisfies (62) for all $k \geq 0$. By continuity of $g(x)$, since $g(x^*) = 0$, there exists $\delta_2 \in (0, \delta_1]$ such that whenever $\|x^k - x^*\| \leq \delta_2$ one has that (61) holds; while, by continuity of $H(x)$, there exists $\delta \in (0, \delta_2]$ such that whenever $\|x - x^*\| \leq \delta$ one has that (60) holds.

On the other hand, by (61), if $\|x^k - x^*\| \leq \delta$, we have that

$$\|g(x^k)\| \leq \left[\frac{1}{2(L + \alpha)} \right] / \|H(x^k)^{-1}\|^2$$

and, since $\|H(x^k)^{-1}\| = 1/\lambda_{1,k}$,

$$\|g(x^k)\| \leq \frac{1}{2(L + \alpha)} \lambda_{1,k}^2.$$

Thus, by Lemma 6.1 and the definition of Algorithm 3.2, we have that x^{k+1} is, in fact, defined by $x^{k+1} = x^k - H_k^{-1}g^k$ and, therefore, the thesis follows by an inductive argument. \square

Theorem 6.2 *Let x^* be a local minimizer of f at which Assumption A3 holds. Suppose also that Assumption A1 holds and, in addition, that x^* is a limit point of the sequence $\{x^k\}_{k=0}^\infty$ generated by Algorithm 3.1–3.2. Then, the whole sequence $\{x^k\}_{k=0}^\infty$ converges quadratically to x^* .*

Proof: Since x^* is a limit point, there exists k_0 such that $\|x^{k_0} - x^*\| \leq \delta$. Thus, the convergence of $\{x^k\}$ follows from Theorem 6.1 replacing x^0 with x^{k_0} . \square

The following is a global non-flatness assumption that will allow us to prove a complexity result that takes advantage of local quadratic convergence.

Assumption A4 *Let $\delta > 0$ be as in the thesis of Theorem 6.1. There exists $\kappa > 0$ such that, for all x^k generated by Algorithm 3.1–3.2, if $\|x^k - x^*\| > \delta$ then $\|g(x^k)\| > \kappa$.*

Note that Assumption A4 holds under the uniform non-singularity assumption that says that for all $k \in \mathbb{N}$ and $x \in [x^k, x^{k+1}]$, $H(x)$ is nonsingular and $\|H(x)^{-1}\| \geq 1/\eta$. In fact, by the Mean Value Theorem, the uniform non-singularity assumption implies that, for all x^k generated by Algorithm 3.1–3.2, $\|g(x^k)\| \geq \eta\|x^k - x^*\|$.

Theorem 6.3 *Let f be bounded below and let x^* be a local minimizer of f at which Assumption A3 holds. Suppose also that Assumptions A1, A2 and A4 hold, and, in addition, that x^* is a limit point of the sequence $\{x^k\}_{k=0}^\infty$ generated by Algorithm 3.1–3.2. Then, after a number of iterations $k_0 = O(\kappa^{-3/2})$, where κ is as in Assumption A4 and it only depends on characteristics of the problem and algorithmic parameters, we obtain that $\|x^k - x^*\| \leq \delta$ for all $k \geq k_0$, where δ is as in the thesis of Theorem 6.1.*

Proof: By construction (see Theorem 6.1), δ only depends on characteristics of the problem. By Assumption A4, $\|g(x^k)\| > \kappa$ for all k such that $\|x^k - x^*\| > \delta$. Then, by Assumptions A1 and A2, and Theorem 5.1, after $k_0 = O(\kappa^{-3/2})$ iterations, we obtain that $\|g(x^{k_0})\| \leq \kappa$, i.e. $\|x^{k_0} - x^*\| \leq \delta$. This implies, by Theorem 6.1, that $\|x^k - x^*\| \leq \delta$ for all $k \geq k_0$, as we wanted to prove. \square

Theorem 6.4 *Let f be bounded below and let x^* be a local minimizer of f at which Assumption A3 holds. Suppose also that Assumptions A1, A2 and A4 hold, and, in addition, that x^* is a limit point of the sequence $\{x^k\}_{k=0}^\infty$ generated by Algorithm 3.1–3.2. Let $\varepsilon_g > 0$ be a given constant. Then, in at most $\hat{k} = O(\log_2(-\log_2(\varepsilon_g)))$ iterations we have that $\|g(x^k)\| \leq \varepsilon_g$ for all $k \geq \hat{k}$.*

Proof: By the Mean Value Theorem of Integral Calculus, we have that, for any $k \geq 0$,

$$g(x^{k+1}) = \left[\int_0^1 H(\xi_{k+1}(t)) dt \right] (x^{k+1} - x^*), \text{ where } \xi_{k+1}(t) = x^* + t(x^{k+1} - x^*). \quad (63)$$

By the triangle inequality and Theorems 6.1 and 6.3, since $\|x^{k+1} - x^*\| \leq \delta$ for all $k \geq k_0$ implies $\|\xi(t) - x^*\| \leq \delta$ for all $k \geq k_0$ and $t \in [0, 1]$, we have that

$$\|H(\xi_{k+1}(t))\| - \|H(x^*)\| \leq \|H(\xi_{k+1}(t)) - H(x^*)\| \leq \gamma \|\xi_{k+1}(t) - x^*\| \leq \gamma\delta \quad (64)$$

for all $k \geq k_0$ and $t \in [0, 1]$. Therefore, by (63) and (64),

$$\|g(x^{k+1})\| = \left\| \left[\int_0^1 H(\xi_{k+1}(t)) dt \right] (x^{k+1} - x^*) \right\| \leq (H(x^*) + \gamma\delta) \|x^{k+1} - x^*\| \quad (65)$$

for all $k \geq k_0$.

On the other hand, by the Mean Value Theorem of Integral Calculus, we have that, for any $k \geq 0$,

$$x^k - x^* = \left[\int_0^1 H(\xi_k(t)) dt \right]^{-1} g(x^k) \text{ where } \xi_k(t) = x^* + t(x^k - x^*)$$

and, thus, by Theorems 6.1 and 6.3, since $\|x^k - x^*\| \leq \delta$ implies $\|\xi_k(t) - x^*\| \leq \delta$ for all $k \geq k_0$ and $t \in [0, 1]$, we have that

$$\|x^k - x^*\| \leq 2\beta \|g(x^k)\| \text{ for all } k \geq k_0. \quad (66)$$

Now, by (65), (66), and Theorems 6.1 and 6.3,

$$\|g(x^{k+1})\| \leq (H(x^*) + \gamma\delta) \|x^{k+1} - x^*\| \leq \beta\gamma(H(x^*) + \gamma\delta) \|x^k - x^*\|^2 \leq 4\beta^3\gamma(H(x^*) + \gamma\delta) \|g^k\|^2 \quad (67)$$

for all $k \geq k_0$.

Up to this point, we have that $\|g^{k_0}\| \leq \kappa$ with $k_0 = O(\kappa^{-3/2})$ and that, for all $\ell \geq 0$, $\|g(x^{k_0+1+\ell})\| \leq c_{\text{quad}}\|g^{k_0+\ell}\|^2$, where κ and $c_{\text{quad}} = 4\beta^3\gamma(H(x^*) + \gamma\delta)$ depend only on characteristics of the problem and algorithmic parameters. This means that

$$\|g(x^{k_0+1+\ell})\| \leq c_{\text{quad}}^{\ell+1}\|g(x^{k_0})\|^{2^{\ell+1}} \leq c_{\text{quad}}^{\ell+1}\kappa^{2^{\ell+1}} \text{ for all } \ell \geq 0. \quad (68)$$

We now consider, with the simple purpose of simplifying the presentation, $k_1 \geq k_0$, $k_1 = O(c_{\text{quad}}^{3/2})$, whose existence is granted by Assumptions A1 and A2 and Theorem 5.1, such that $\|g^k\| \leq \frac{1}{2}c_{\text{quad}}^{-1}$ for all $k \geq k_1$. Thus, (68) can be re-stated as

$$\|g(x^{k_1+1+\ell})\| \leq c_{\text{quad}}^{\ell+1}\|g(x^{k_1})\|^{2^{\ell+1}} \leq \frac{c_{\text{quad}}^{\ell+1}}{c_{\text{quad}}^{2^{\ell+1}}} \left(\frac{1}{2}\right)^{2^{\ell+1}} \leq 2^{-2^{\ell+1}} \text{ for all } \ell \geq 0. \quad (69)$$

Thus, since $2^{-2^{\ell+1}} \leq \varepsilon_g$ if and only if $\ell \geq \log_2(-\log_2(\varepsilon_g)) + 1$, we have that $\|g^k\| \leq \varepsilon_g$ for all $k \geq k_1 + \log_2(-\log_2(\varepsilon_g)) + 1$. This implies the desired result recalling that k_1 does not depend on ε_g . \square

7 Numerical experiments

We implemented Algorithm 3.1–3.2 in Fortran 90. At each iteration k , the spectral decomposition of matrix H_k is computed by the Lapack [1] subroutine DSYEV. At Step 5.1 of Algorithm 3.2, $\tilde{\mu}_{k,\ell_5} > 0$ and \tilde{z}^{k,ℓ_5} solution to (10) with $\mu = \tilde{\mu}_{k,\ell_5}$ such that (14) holds are computed using bisection. In the numerical experiments, we arbitrary considered $\alpha = 10^{-8}$ and $M = 10^3$. It should be noted that these two parameters, as well as the other constants that appeared hard-coded in Algorithm 3.1–3.2 (in order to simplify the exposition), were not subject to tuning at all. All those values were chosen because they seemed to be “natural choices” and the intention of the numerical experiments below is not to deliver the most robust or efficient version of the proposed method but to illustrate its practical behaviour.

The method proposed in the present work will be compared against the line-search Newton’s method with quadratic regularization and Armijo descent introduced in [16]. With this purpose, we implemented (also in Fortran 90) Algorithm 1 described in [16, p.348]. In order to focus the comparison on the methods’ differences (mainly the way in which the regularizing parameter is computed and the descent criterion), our implementation uses the Lapack subroutine DSYEV for computing the spectral decomposition of H_k . This choice provides the value of the smallest eigenvalue of H_k required by the algorithm and also trivializes solving the Newtonian linear system. A classical quadratic interpolation (taking $t/2$ as a new trial step when the minimizer of the quadratic model lies outside the interval $[0.1t, 0.9t]$) was considered. In the numerical experiments, we set, as suggested in [16], $\beta = 10^{-2}$, $\eta = 0.25$, $L_0 = 10^{-6}$, and $\delta = 10^{-16}$. We considered the two choices $\mu_k = \mu_k^-$ and $\mu_k = \mu_k^+$ and, thus the method introduced in [16] with these two choices will be referred, from now on, as “KKS with $\mu_k = \mu_k^-$ ” and “KKS with $\mu_k = \mu_k^+$ ”.

The Fortan 90 implementation of Algorithm 3.1–3.2, as well as our implementation of the algorithm introduced in [16], is freely available at <http://www.ime.usp.br/~egbrigin/>. Interfaces for solving user-defined problems coded in Fortran 90 as well as problem from the CUTEst

collection [13] are available. All tests reported below were conducted on a computer with 3.5 GHz Intel Core i7 processor and 16GB 1600 MHz DDR3 RAM memory, running OS X Yosemite (version 10.10.5). Codes were compiled by the GFortran compiler of GCC (version 5.1.0) with the `-O3` optimization directive enabled.

7.1 An *ad hoc* toy problem with expected hard case

In this section, we illustrate the behaviour of Algorithm 3.1–3.2 in a simple problem in which the hard case is expected to appear. Consider the function defined by

$$f(x_1, x_2) = x_1x_2 + 0.1(x_1 - x_2)^4 + (x_1 + x_2)^4, \quad (70)$$

whose level sets are displayed in Figure 2. This function has two global minimizers at, approximately, $(0.559017, -0.559017)$ and $(-0.559017, 0.559017)$, at which the functional value is approximately -0.15625 . Moreover, $(0, 0)$ is a saddle point at which f vanishes. We are interested in the behaviour of the considered algorithms when the initial point is in the line $x_1 = x_2$ and relatively close to $(0, 0)$.

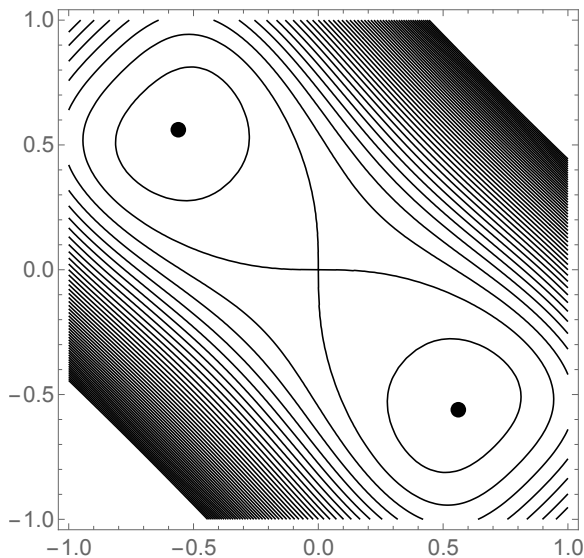


Figure 2: Contour plot of $f(x_1, x_2)$ defined in (70).

The Hessian is indefinite if $x_1 = x_2$ and the eigenvalues of $\nabla^2 f(x_1, x_2)$ tend to 1 and -1 when $x_1 = x_2$ and $x_1 \rightarrow 0$. For all iterate satisfying $x_1 = x_2$ the minimum norm solution of (10) satisfies $s_1 = s_2 \approx -x_1 = -x_2$. Since the regularization parameter tends to 1 when $x_1 = x_2$ and $x_1 \rightarrow 0$, it turns out that the associated ρ tends to infinity when $x_1 = x_2$ and $x_1 \rightarrow 0$. As a consequence, when an iterate (x_1^k, x_2^k) with $x_1^k = x_2^k$ is close to the origin, the test $\rho_{k,0} \leq M$ eventually fails at Step 2 of Algorithm 3.2 and a search along the eigenvector orthogonal to $x_1 = x_2$ is initiated. So, the process quickly converges to one of the global minimizers. On the other hand, a Newtonian method like the one considered in [16] never leaves the line $x_1 = x_2$ and convergence to the saddle point $(0, 0)$ is expected.

If we run Algorithm 3.1–3.2 starting from $(x_1^0, x_2^0) = (1, 1)$, for all iterations $k \leq 14$, we observe that, in fact, the linear system (10) is compatible, $\rho_{k,0} \leq M$, and $\hat{s}^{k,0}$ satisfies the descent condition (8). Therefore, we have that $x^{14} \approx (2.53523, 2.53523) \times 10^{-4}$ still lies in the line $x_1 = x_2$. At iteration $k = 15$, we have that $\rho_{k,0} > M$ and a search along the eigenvector is performed. Having abandoned the line $x_1 = x_2$, convergence to the global minimizer $(-0.559017, 0.559017)$ occurs and the algorithm stops at iteration $k = 20$ satisfying $\|\nabla f(x^{20})\|_\infty \leq 10^{-8}$ and $\lambda_1(\nabla^2 f(x^{20})) \geq -10^{-8}$ and performing, as a whole, 23 functional evaluations and having solved 30 linear systems.

Methods KKS with $\mu_k = \mu_k^-$ and KKS with $\mu_k = \mu_k^+$, as expected, converge to the saddle point $(0, 0)$ (using only two iterations, three functional evaluations, and solving three linear systems). The considered *ad hoc* problem was presented in order to highlight a property of the proposed method (related to robustness) that may not be shared by other methods. Since different final iterates are being found, it would be meaningless to compare the effort required by each method for achieving a stopping criterion (first- or second-order criticality); while ignoring the objective functional value at the final iterate.

If we now run Algorithm 3.1–3.2 starting from $(0, 0)$, it converges to the same global minimizer in 9 iterations using 11 functional evaluations and having solved 18 linear systems; while, as expected, methods KKS with $\mu_k = \mu_k^-$ and KKS with $\mu_k = \mu_k^+$ satisfy the stopping criteria at the initial point.

7.2 Massive comparison

In this section we consider the 87 problems from the CUTEst collection already considered in the numerical experiments presented in [16]. The same dimensions chosen in [16] were preserved (most of the problems have $n = 1000$ variables). These problems correspond to *all* the unconstrained problems from the CUTEst collection with available second-order derivatives.

For the stopping criteria, we set $f_{\min} = -10^{10}$, $\varepsilon_g^a = 10^{-6}$, and $\varepsilon_g^r = 10^{-15}$. Other than stopping if an iterate x^k satisfies

$$f(x^k) \leq f_{\min}$$

or

$$\|g^k\| \leq \varepsilon_g^a, \tag{71}$$

the methods also stop if

$$\|g^k\| \leq \varepsilon_g^r \|g^0\| \tag{72}$$

or if the elapsed CPU time exceeds one hour. It should be noted that, in order to allow a fair comparison, the same first-order criticality stopping criteria are being used for KKS with $\mu_k = \mu_k^-$ and KKS with $\mu_k = \mu_k^+$ as well as for Algorithm 3.1–3.2. However, this choice does not affect the quality of the final points obtained by Algorithm 3.1–3.2 because a simple inspection of the results reveals that, in the considered set of problems, any time the stopping criteria (71) or (72) is satisfied, its second-order counterpart, given by

$$\|g^k\| \leq \varepsilon_g^a \text{ and } \lambda_{1,k} \geq -\varepsilon_H^a$$

and

$$\|g^k\| \leq \varepsilon_g^r \|g^0\| \text{ and } \lambda_{1,k} \geq -\varepsilon_H^r \max_{j=1,n} \{|\lambda_{j,0}|\}$$

(with $\varepsilon_H^a = \varepsilon_g^a$ and $\varepsilon_H^r = \varepsilon_g^r$), respectively, is satisfied as well. We will refer to these stopping criteria as 'UN' (unbounded f), 'AS' (first- or second-order absolute stopping), 'RS' (first- or second-order relative stopping), and 'TE' (CPU time limit exceeded). Exceptionally, although $\|\cdot\|$ stands for the Euclidean norm everywhere in the text, the sup-norm of the gradient was considered at the stopping criteria described above. None other stopping criterion was considered.

Detailed information regarding the performance of each method on each problem can be found at <http://www.ime.usp.br/~egbrigin/>. For a given problem, let f_1 , f_2 , and f_3 be the value of the objective function at the final iterate delivered by each of the three methods. Following [3], we will say that the three methods found *equivalent solutions* if

$$\frac{f_i - f_{\text{best}}}{\max\{1, |f_{\text{best}}|\}} \leq 10^{-2} \text{ for } i = 1, 2, 3,$$

where $f_{\text{best}} = \min\{f_1, f_2, f_3\}$. The 87 problems will be separated into two sets. Set 1 will be given by the 66 problems in which the three methods found equivalent solutions and stopped satisfying the absolute or the relative stopping criterion. Set 2 will contain the remaining 21 problems. Problems in Set 1 will be used to analyze the efficiency of the methods; while problems in Set 2 will be observed with an eye on robustness.

For analyzing the efficiency of the methods through its performance on the 66 problems on Set 1, we used performance profiles [10]. See Figure 3. By definition of the performance profiles and the way in which the problems were selected, all curves reach the value 1 at the right-hand-side of the graphic. Thus, these pictures evaluate efficiency only. The three pictures show the same thing: Algorithm 3.1–3.2 is more efficient in most of the problems but there are a few problems in which it takes much longer than the other two methods.

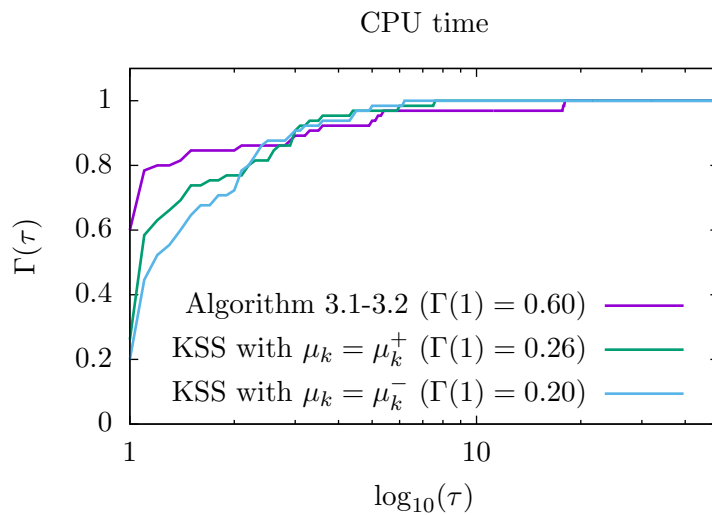
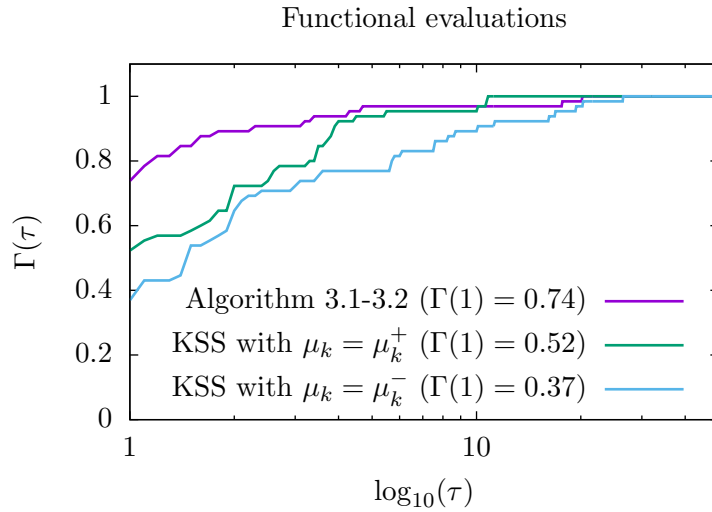
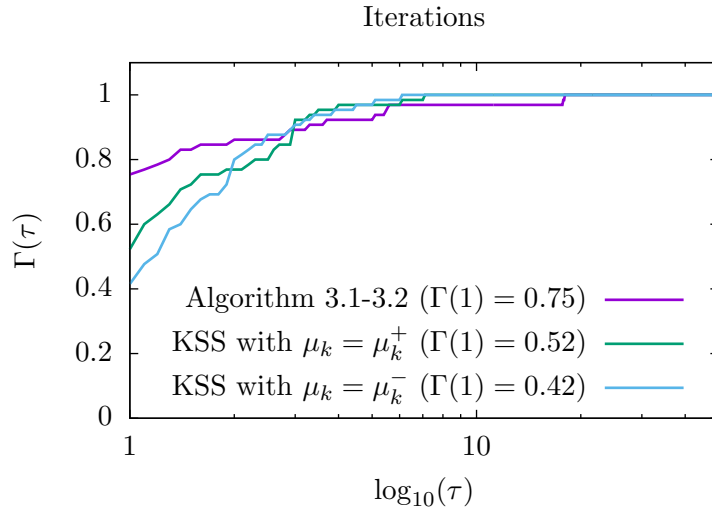


Figure 3: Performance profiles considering the 66 problems in which the three methods stopped satisfying the same stopping criterion related to absolute or relative criticality and found equivalent solutions.

Table 1 shows the details of the final iterates found by the three methods on problems in Set 2. It can be said that, considering these 21 problems, Algorithm 3.1–3.2 satisfied the second-order criticality stopping criteria 13 times; while KSS with $\mu_k = \mu_k^-$ and KSS with $\mu_k = \mu_k^+$ satisfied the first-order criticality stopping criteria 5 and 11 times, respectively. Other than that, there are 3 problems (FLETGBV3, FLETCHBV, INDEF) in which the objective function appears to be unbounded from below. KSS with $\mu_k = \mu_k^-$ and KSS with $\mu_k = \mu_k^+$ were both able to identify this situation and stopped by the UN stopping criterion. Algorithm 3.1–3.2 recognized the situation in only one of the cases and stopped by TE in the other two. This may indicate that Algorithm 3.1–3.2 takes longer to reduce the objective functional value when it is unbounded below. There are also cases in which the three methods found an approximate stationary point but did not find equivalent solutions. BROYDN7D, CHAINWOO, and NCB20 are examples of these cases. The methods take turn to be the one that finds the stationary point with the lowest functional value and, therefore, the presented experiment did not show whether any of the methods is able to find better quality solutions.

Problem name	Algorithm 3.1–3.2			KSS with $\mu_k = \mu_k^-$			KSS with $\mu_k = \mu_k^+$		
	$f(x^k)$	$\ g^k\ $	SC	$f(x^k)$	$\ g^k\ $	SC	$f(x^k)$	$\ g^k\ $	SC
BROYDN7D	3.54624D+02	2.1D-10	AS	4.81627D+02	1.6D-11	AS	4.60601D+02	6.7D-07	AS
CHAINWOO	1.57548D+02	2.1D-12	AS	1.00000D+00	1.7D-12	AS	1.00000D+00	1.5D-09	AS
COSINE	-9.99000D+02	1.1D-12	AS	-1.40035D+02	2.4D+04	TE	-9.44546D+02	1.6D+00	TE
ENGVAl1	1.10819D+03	1.3D-12	AS	1.10819D+03	1.3D-12	AS	1.10819D+03	1.8D-06	TE
FLETGBV3	-1.54153D+03	3.0D-02	TE	-1.00026D+08	1.2D-01	UN	-1.00026D+08	1.4D-01	UN
FLETCHBV	-1.84122D+09	2.8D+06	UN	-1.84122D+09	2.8D+06	UN	-1.84122D+09	2.8D+06	UN
GENHUMPS	8.73814D+06	1.1D+02	TE	5.90238D+06	1.3D+02	TE	7.70165D+06	1.5D+02	TE
INDEF	-2.72320D+06	1.0D+00	TE	-1.09591D+08	1.0D+00	UN	-1.09760D+08	1.0D+00	UN
MANCINO	1.67148D-14	1.0D-03	RS	2.14315D+17	3.0D+12	TE	1.67797D-14	5.5D-04	RS
MODBEALE	1.10832D-20	9.5D-10	AS	5.19223D+01	1.8D-04	TE	8.04120D+00	1.7D-05	TE
NCB20	9.32122D+02	4.5D-10	AS	9.16688D+02	5.9D-07	AS	9.17763D+02	5.6D-08	AS
NONCVXUN	2.32878D+03	1.6D-03	TE	2.32595D+03	3.4D-08	AS	2.31974D+03	1.4D-07	AS
NONMSQRT	9.02177D+01	3.6D-04	TE	8.99049D+01	3.1D-01	TE	8.99048D+01	4.4D-01	TE
PENALTY2	1.12970D+83	3.4D+75	TE	1.44640D+83	2.1D+38	TE	1.44640D+83	2.1D+38	TE
PENALTY3	9.99523D-04	1.2D-07	AS	3.98575D+04	8.7D-02	TE	9.94993D-04	7.2D-04	TE
SBRYBND	8.80296D-27	3.5D-06	TE	2.49040D+04	2.0D+07	TE	1.85974D-21	6.8D-07	AS
SCOSINE	1.09888D+02	2.9D+13	TE	8.76705D+02	1.2D+05	TE	8.57518D+02	1.2D+11	TE
SCURLY10	-1.00316D+05	4.3D-08	AS	0.00000D+00	1.8D+05	TE	-1.00316D+05	1.5D-07	AS
SCURLY20	-1.00316D+05	1.4D-07	AS	0.00000D+00	3.4D+05	TE	-1.00316D+05	1.2D-07	AS
SCURLY30	-1.00316D+05	1.1D-07	AS	0.00000D+00	5.0D+05	TE	-1.00316D+05	3.1D-07	AS
SENSORS	-2.10853D+05	6.8D-10	AS	-2.10916D+05	1.7D-05	TE	-2.10633D+05	1.1D-09	AS
SPMSRTLS	4.34760D-16	3.2D-11	AS	4.37365D-16	3.1D-09	AS	1.75675D+00	2.4D-07	AS

Table 1: Details of the 21 problems in which it does not hold that “the three methods stopped satisfying the first- or second-order criticality stopping criterion and found equivalent solutions”.

8 Final remarks

Iterative regularization is a traditional idea for unconstrained optimization. Levenberg [17] and Marquardt [18] were the first to apply this idea to nonlinear least-squares problems with the main purpose of stabilizing the Gauss-Newton method. In the nonlinear least-squares context the Hessian of a natural quadratic approximation of the objective function has the form $J(x)^T J(x)$, where $J(x)$ is the Jacobian of the vectorial function whose squared norm one wants to minimize.

This matrix is positive semidefinite but may be very ill-conditioned if the independent variables of the underlying approximation problem are highly correlated. Therefore, replacing $J(x)^T J(x)$ with $J(x)^T J(x) + \mu I$ has the beneficial effect of adding stability to the quadratic model and, perhaps, also to the original minimization problem.

When the quadratic model is built with the full Hessian $H(x)$ of the objective function we face the additional difficulty that this matrix may be indefinite, so that minimizers of the quadratic approximation may not exist at all. In this case, one may replace $H(x)$ with $H(x) + \sigma I$ in the quadratic model, where $\sigma \geq \max\{0, -\lambda_1(H(x))\}$ and $\lambda_1(H(x))$ denotes the smallest eigenvalue of the matrix. If $\lambda_1(H(x)) > 0$ or $\sigma > \max\{0, -\lambda_1(H(x))\}$, the quadratic model has exactly one solution. On the other hand, if $\lambda_1(x) \leq 0$ and $\sigma = -\lambda_1(x)$ two situations may occur: either the regularized quadratic model has no solution or it has infinitely many solutions. The second case characterizes the hard case, in which the set of minimizers of the regularized model for different values of σ has two branches. Quadratic regularization is attractive because, given a regularizing parameter $\sigma > \max\{0, -\lambda_1(H(x))\}$ or if $\lambda_1(H(x)) > 0$, the only minimizer of the model that provides a suitable trial point comes from the solution of a single linear system. However, in the hard case, the minimum norm solution of the model could be excessively small. For this reason, in this case, it is recommendable to consider the existent infinitely many solutions in order to produce a more aggressive exploration of the space with the aim of catching global minimizers.

It is straightforward to verify that any minimizer \bar{s} of the regularized model for some parameter σ is also a solution of the non-regularized model on the region $\|s\| \leq \|\bar{s}\|$. This property characterizes the essential equivalence between trust-region and regularization ideas. The equivalence extends to the case of cubic regularization and, in turn, gives rise to optimal complexity results.

The present paper explored all these ideas with the principal objective of preserving the complexity results that hold in the case of cubic regularization. Although there are good algorithms for solving the cubic regularization subproblem, these algorithms, as well as the ones for solving the trust-region subproblem, generally need to solve more than one linear system for computing a trial point. Unfortunately, in the algorithm introduced in this paper we could not preserve the property of “one linear system per trial point” at every iteration, because the preservation of complexity needed safeguarded choices for computing the first nonnull regularization parameter μ . On the other hand, even a preliminary implementation in which algorithmic parameters were not tuned at all, produced satisfactory results, in comparison with a well-established regularization method for unconstrained optimization. In addition to first- and second-order complexity results, we proved asymptotic convergence to first- and second-order stationary points, as well as local convergence and a complexity result corresponding to the case in which local quadratic convergence takes place.

The regularization method introduced in [16] and our present regularized method were conceived with quite different purposes. While in our case we were worried about the compatibility of the most simple updating rules of the regularization parameter with the preservation of optimal complexity results, in [16] the main concern was the determination of regularizing parameters that optimize the accuracy of the quadratic model. The natural challenge that emerges is related, therefore, with the compatibility between the updating rules of [16] and our updating rules and purposes. It should be mentioned, moreover, that in [16] a line search follows the obtention of the adequate point on the Levenberg-Marquardt path, motivating additional ques-

tions about the compatibility of this search with complexity bounds. Needless to say, this type of studies should be complemented with insightful and extensive numerical experiments.

References

- [1] E. Anderson, Z. Bai, C. Bischof, S. Blackford, J. Demmel, J. Dongarra, J. Du Croz, A. Greenbaun, S. Hammarling, A. McKenney, and D. Sorensen, *LAPACK Users' Guide*, third edition, Society for Industrial and Applied Mathematics, Philadelphia, 1999.
- [2] E. G. Birgin, J. L. Gardenghi, J. M. Martínez, S. A. Santos, and Ph. L. Toint, Worst-case evaluation complexity for unconstrained nonlinear optimization using high-order regularized models, *Mathematical Programming*, to appear (DOI: 10.1007/s10107-016-1065-8).
- [3] E. G. Birgin and J. M. Gentil, Evaluating bound-constrained minimization software, *Computational Optimization and Applications* 53, pp. 347–373, 2012.
- [4] 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, pp. 2833–2852, 2010.
- [5] C. Cartis, N. I. M. Gould, and Ph. L. Toint, Adaptive cubic regularization methods for unconstrained optimization. Part I: motivation motivation, convergence and numerical results, *Mathematical Programming* 127, pp. 245–295, 2011.
- [6] C. Cartis, N. I. M. Gould, and Ph. L. Toint, Adaptive cubic regularization methods for unconstrained optimization. Part II: worst-case function and derivative complexity, *Mathematical Programming* 130, pp. 295–319, 2011.
- [7] A. R. Conn, N. I. M. Gould, and Ph. L. Toint, *Trust Region Methods*, Society for Industrial and Applied Mathematics, Philadelphia, 2000.
- [8] F. E. Curtis, D. P. Robinson, and M. Samadi, A trust-region algorithm with a worst-case iteration complexity of $O(\varepsilon^{-3/2})$, *Mathematical Programming*, to appear (DOI: 10.1007/s10107-016-1026-2).
- [9] J. E. Dennis Jr. and R. B. Schnabel, *Numerical Methods for Unconstrained Optimization and Nonlinear Equations*, Society for Industrial and Applied Mathematics, Philadelphia, 1996.
- [10] E. D. Dolan and J. J. Moré, Benchmarking optimization software with performance profiles, *Mathematical Programming* 91, pp. 201–213, 2002.
- [11] J. P. Dussault, Simple unified convergence proofs for the trust-region and a new ARC variant, Technical report, University of Sherbrooke, Sherbrooke, Canada, 2015.
- [12] C. C. Gonzaga and E. W. Karas, Complexity of first-order methods for differentiable convex optimization, *Pesquisa Operacional* 34, pp. 395–419, 2014.

- [13] N. I. M. Gould, D. Orban, and Ph. L. Toint, CUTEst: a constrained and unconstrained testing environment with safe threads for mathematical optimization, *Computational Optimization and Applications* 60, pp. 545–557, 2014.
- [14] G. N. Grapiglia, J-Y Yuan, and Y-X Yuan, On the convergence and worst-case complexity of trust-region and regularization methods for unconstrained optimization, *Mathematical Programming* 152, pp. 491–520, 2015.
- [15] A. Griewank, *The modification of Newton’s method for unconstrained optimization by bounding cubic terms*, Technical Report NA/12, Department of Applied Mathematics and Theoretical Physics, University of Cambridge, 1981.
- [16] E. W. Karas, S. A. Santos, and B. F. Svaiter, Algebraic rules for quadratic regularization of Newton’s method, *Computational Optimization and Applications* 60, pp. 343–376, 2015.
- [17] K. Levenberg, A method for the solution of certain non-linear problems in least-squares, *Quarterly Journal of Applied Mathematics* 2, pp. 164–168, 1944.
- [18] D. Marquardt, An algorithm for least-squares estimation of nonlinear parameters, *SIAM Journal on Applied Mathematics* 11, pp. 431–441, 1963.
- [19] J. M. Martínez and M. Raydan, Cubic-regularization counterpart of a variable-norm trust-region method for unconstrained minimization, *Journal of Global Optimization*, to appear (DOI: 10.1007/s10898-016-0475-8).
- [20] Y. Nesterov and B. T. Polyak, Cubic regularization of Newton’s method and its global performance, *Mathematical Programming* 108, pp. 177–205, 2006.