





Error estimates for iterative algorithms for minimizing regularized quadratic subproblems

Nicholas I. M. Gould ^a and Valeria Simoncini ^{b,c}

^aScientific Computing Department, STFC-Rutherford Appleton Laboratory, Chilton, England; ^bDipartimento di Matematica, Università di Bologna, Bologna, Italy; ^cIMATI-CNR, Pavia, Italy

ABSTRACT

We derive bounds for the objective errors and gradient residuals when finding approximations to the solution of common regularized quadratic optimization problems within evolving Krylov spaces. These provide upper bounds on the number of iterations required to achieve a given stated accuracy. We illustrate the quality of our bounds on given test examples.

ARTICLE HISTORY

Received 4 April 2019
Accepted 17 September 2019

KEYWORDS

Trust-region subproblem; regularized quadratic subproblem; error estimates; Krylov subspace

2010 MATHEMATICS SUBJECT CLASSIFICATION
90C30

1. Introduction

In this paper, we derive upper bounds for the number of iterations required to reach a certain level of optimality by subspace methods for solving the trust-region subproblem

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad q(x) := g^T x + \frac{1}{2} x^T H x \quad \text{subject to } \|x\| \leq \delta \quad (1)$$

and its regularization variant

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad q^R(x, \sigma, p) := q(x) + \frac{1}{p} \sigma \|x\|^p. \quad (2)$$

Here, we are given a gradient g , a symmetric, but possibly indefinite, Hessian H , a radius $\delta > 0$, a weight $\sigma > 0$ and a power $p > 2$, and use the Euclidean norm $\|\cdot\|$. Subproblems (1)–(2) lie at the heart of the step calculation in both trust-region and cubic-regularization methods for unconstrained optimization [9,11,31,32].

A typical requirement in the trust-region case is that the computed x should decrease the objective function, i.e. $q(x) < q(0) \equiv 0$, and that the gradient of the Lagrangian for the problem, $g + Hx + \mu x$, should be smaller than a prescribed tolerance in norm, i.e.

$$\|g + Hx + \mu x\| \leq \epsilon \quad (3)$$

for some given $\epsilon > 0$, whose precise value determines the rate of convergence of the trust-region algorithm, and a suitable Lagrange multiplier, $\mu \geq 0$, for the trust-region constraint $\|x\| \leq \delta$. For regularization problems, a similar requirement is that $q^R(x, \sigma, p) < q^R(0, \sigma, p) \equiv 0$ and that the norm of the gradient of $q^R(x, \sigma, p)$ should be small. Since $\nabla_x q^R(x, \sigma, p) = g + Hx + \mu x$ where $\mu = \sigma \|x\|^{p-2}$, the latter requirement is identical to (3) but for a different μ . As the subspace methods we consider automatically ensure that their relevant objectives decrease, our intention is to provide bounds on the number of steps (actually products with H) required by such methods to achieve (3) for the problems under consideration.

The subspaces of interest here are the nested Krylov spaces $\mathcal{K}_k := \mathcal{K}(H, g, k)$ for $k \geq 0$, where, for general A and b , we define $\mathcal{K}(A, b, k) := \text{span}\{A^i b\}_{i=0}^{k-1}$. A sequence of estimates x_k are generated so that

$$x_k = \arg \min_{x \in \mathcal{K}_k} q(x) \quad \text{subject to } \|x\| \leq \delta \quad (4)$$

for the trust-region subproblem, or

$$x_k = \arg \min_{x \in \mathcal{K}_k} q^R(x, \sigma, p) \quad (5)$$

for the regularization case; here and elsewhere $\arg \min$ refers to the set of global minimizers of the function under consideration within the domain specified.¹ This is useful as the well-known GLTR method [18] for (1) and the GLRT approach [9] for (2), which exploit the evolving Lanczos basis for \mathcal{K}_k , use precisely these formulations.² However, we must be cautious as it is well known [18, Theorem 5.8] that such methods may fail to solve the problem if the sequence of Krylov subspaces lies in an unpropitious non-trivial invariant subspace of \mathfrak{R}^n , and in this case it may be necessary to enhance the search space with a specific eigenvector of H from outside the Krylov space. Fortunately, as we shall see, this is not necessary if our goal is merely to satisfy (3).

To put subspace methods in context, we briefly review other approaches. Early methods [16,30] were aimed at the trust-region case (1), and use the optimality conditions for the problem (see Theorem 2.1) to reduce it to a scalar root-finding problem involving a so-called secular equation. A safeguarded Newton process is applied, and the required value and derivatives of the related secular function at each scalar value μ require the direct solution of a symmetric, positive-definite linear system $(H + \mu I)x = -g$. These ideas may be extended to higher-order root-finding methods, and for problems with sparse Hessians [20].

The first methods designed to avoid factorization [37,38] are only able to find approximate solutions, based on properties of the conjugate-gradient method. GLTR [18] is the first method fully to explore this Krylov subspace. Other subspace methods [13,14,23] have focussed on evolving very low-dimensional (non-Krylov) subspaces. Another popular approach is to couch (1) as a parametric or generalized eigenvalue problem [1,26,34,35], and to harness existing, powerful software for this domain. A third suggestion is to formulate the problem as a semi-definite or second-order cone program [15,24,25,33], and as before to use sophisticated methods from this area. GLTR has a good reputation [1] unless the solution required lies outside the Krylov space, when, perhaps unsurprisingly, eigenvalue approaches may have the edge.

The regularization subproblem (5) has risen in importance in the 21st century, and obvious analogues [6,9,20] of the direct and subspace approaches have been proposed. There has also been a move towards simpler (accelerated, first-order methods) [5,7,40] that recognize that matrices may be too large to manipulate when the problem is enormous.

While all of these methods provide guarantees of convergence to a solution of their relevant subproblem, little has been said about how this convergence occurs. Recently there has been some effort to bound the decrease in the objective function value, and thus the overall number of iterations required to reach a specified objective accuracy, for first-order and subspace methods [8,24,25,41,42]. Our interest, though, is in trying to achieve (3), since this is more significant for the overall convergence of many nonlinear optimization methods [9,11,32]. This is the main focus of our paper.

In §2 we examine the benefits and limitations of Krylov approximations to the solutions we wish to find. We follow this, in §3, by deriving bounds both on the decrease in the model objective functions and on the norm of the violation of the first-order criticality residuals from the Krylov space under consideration. We note that bounds on the former have already been proposed [8,41,42], and indeed our bounds make use of some of those in [8]. The bounds we derive for the residuals generalize well-known ones for conjugate-gradient methods applied to definite linear systems. They behave just as in the conjugate-gradient case, but necessarily reflect the additional complications due to the potential indefiniteness of the matrices involved and the presence of explicit or implicit regularization. The bounds allow us to give an upper estimate of how many iterations will be required to satisfy (3), and thus on the work involved in finding a suitable approximation to the solution of the subproblem under consideration. We examine our residual bounds on test examples that are designed to illustrate a variety of spectral distributions in §4. Finally, we make concluding remarks in §5.

2. Solutions from the Krylov space and beyond

Let $\lambda_1 \leq \dots \leq \lambda_n$ be the eigenvalues of H , with the leftmost λ_1 having multiplicity n_1 , and let $u_i, i \in \mathcal{N} := \{1, \dots, n\}$ be the corresponding orthonormal eigenvectors. Crucially, there are well-known characterizations of the global solutions of (1) and (2).

Theorem 2.1 ([16], **Theorem 2.1**; [30], **Lemma 2.1**): *Any global solution x_* to the trust-region subproblem (1) satisfies*

$$(H + \mu_* I)x_* = -g, \tag{6}$$

where the Lagrange multiplier $\mu_* \geq \max(0, -\lambda_1)$ and $\mu_*(\|x_*\|^2 - \delta^2) = 0$. Moreover the solution is unique and $\mu_* > \max(0, -\lambda_1)$ whenever $g^T u_i \neq 0$ for some $1 \leq i \leq n_1$.

Theorem 2.2 ([9], **Theorem 3.1**; [31], **Theorem 10**): *Any global solution x_* to the regularization subproblem (2) satisfies (6), where the multiplier $\mu_* = \sigma \|x_*\|^{p-2} \geq -\lambda_1$. Moreover the solution is unique and $\mu_* > -\lambda_1$ whenever $g^T u_i \neq 0$ for some $1 \leq i \leq n_1$.*

We consider the evolving Krylov spaces \mathcal{K}_k , $k \geq 0$, in more detail. Clearly we may decompose

$$g = \sum_{j=1}^n (g^T u_j) u_j$$

in terms of the basis of eigenvectors $\{u_j\}_{j \in \mathcal{N}}$ of H . Let $\mathcal{I}_+ := \{j \mid g^T u_j \neq 0\}$, $\mathcal{I}_0 := \mathcal{N} \setminus \mathcal{I}_+$ and $m := |\mathcal{I}_+|$.³ Thus

$$g = \sum_{j \in \mathcal{I}_+} (g^T u_j) u_j \quad \text{and hence} \quad H^i g = \sum_{j \in \mathcal{I}_+} \lambda_j^i (g^T u_j) u_j.$$

Therefore $\mathcal{K}_m = \dots = \mathcal{K}_n$, since $\mathcal{K}_m = \text{span}\{u_j\}_{j \in \mathcal{I}_+}$ and the vectors $H^i g$ for $m < i \leq n$ are dependent on those in \mathcal{K}_m . Hence, our Krylov methods will make no further progress beyond the m th iteration, and at that point provide estimates of their relevant solutions x_m and multipliers μ_m .

We now contrast x_m with a desired global minimizer x_* . Let U_+ be the n by m matrix whose columns are the eigenvectors u_j , $j \in \mathcal{I}_+$, U_0 be the n by $n-m$ matrix whose columns are the remaining eigenvectors and $U = (U_+ : U_0)$. Likewise let Λ be the diagonal matrix of eigenvalues ordered as for U , and let Λ_+ and Λ_0 be its diagonal blocks. Thus

$$\Lambda = U^T H U = \begin{pmatrix} \Lambda_+ & 0 \\ 0 & \Lambda_0 \end{pmatrix}. \quad (7)$$

If we define $\bar{g} := U^T g$, and therefore $g = U \bar{g}$, this leads to

$$\begin{pmatrix} \bar{g}_+ \\ \bar{g}_0 \end{pmatrix} := \bar{g} = \begin{pmatrix} U_+^T g \\ U_0^T g \end{pmatrix} = \begin{pmatrix} U_+^T g \\ 0 \end{pmatrix} \quad \text{and} \quad g = U_+ \bar{g}_+, \quad (8)$$

since $\bar{g}_0 = 0$ as $u_j^T g = 0$ for all $j \in \mathcal{I}_0$.

Consider the trust-region subproblem (1), and the change of variables $x = U \bar{x}$. In this case, $x_* = U \bar{x}_*$, where

$$\bar{x}_* \in \arg \min_{\bar{x} \in \mathbb{R}^n} \bar{g}^T \bar{x} + \frac{1}{2} \bar{x}^T \Lambda \bar{x} \quad \text{subject to} \quad \|\bar{x}\| \leq \delta. \quad (9)$$

The optimality conditions (6) for this are

$$\begin{pmatrix} \Lambda_+ + \mu_* I & 0 \\ 0 & \Lambda_0 + \mu_* I \end{pmatrix} \begin{pmatrix} \bar{x}_*^+ \\ \bar{x}_*^0 \end{pmatrix} = - \begin{pmatrix} \bar{g}_+ \\ 0 \end{pmatrix}, \quad (10)$$

where \bar{x}_* and the Lagrange multiplier

$$\mu_* \geq \max(0, -\lambda_1) \quad (11)$$

satisfy

$$\|\bar{x}_*\| \leq \delta \quad \text{and} \quad \mu_* (\|\bar{x}_*\|^2 - \delta^2) = 0, \quad (12)$$

and we have partitioned

$$\bar{x}_* = U^T x_* \equiv \begin{pmatrix} \bar{x}_*^+ \\ \bar{x}_*^0 \end{pmatrix}.$$

By contrast, if $x \in \mathcal{K}_m$, then $x = U_+ \hat{x}_+$ for some vector $\hat{x}_+ \in \mathbb{R}^m$, in which case (4) gives $x_m = U_+ \hat{x}_+^+$, where uniquely

$$\hat{x}_*^+ = \arg \min_{\hat{x}_+ \in \mathbb{R}^m} \bar{g}_+^T \hat{x}_+ + \frac{1}{2} \hat{x}_+^T \Lambda_+ \hat{x}_+ \quad \text{subject to } \|\hat{x}_+^+\| \leq \delta. \quad (13)$$

The optimality conditions (6) then imply that

$$(\Lambda_+ + \mu_*^+ I) \hat{x}_*^+ = -\bar{g}_+, \quad (14)$$

where \hat{x}_*^+ and the Lagrange multiplier

$$\mu_m \equiv \mu_*^+ > \max \left(0, -\min_{j \in \mathcal{I}_+} \lambda_j \right) \quad (15)$$

satisfy

$$\|\hat{x}_*^+\| \leq \delta \quad \text{and} \quad \mu_*^+ (\|\hat{x}_*^+\|^2 - \delta^2) = 0. \quad (16)$$

Given \hat{x}_*^+ , let $\hat{x}_*^0 = 0$ and define

$$\hat{x}_* = \begin{pmatrix} \hat{x}_*^+ \\ \hat{x}_*^0 \end{pmatrix} \quad \text{so that } x_m = U \hat{x}_*.$$

In this case, (14) and (16) become

$$\begin{pmatrix} \Lambda_+ + \mu_*^+ I & 0 \\ 0 & \Lambda_0 + \mu_*^+ I \end{pmatrix} \begin{pmatrix} \hat{x}_*^+ \\ \hat{x}_*^0 \end{pmatrix} = -\begin{pmatrix} \bar{g}_+ \\ 0 \end{pmatrix}, \quad (17)$$

and

$$\|\hat{x}_*\| \leq \delta \quad \text{and} \quad \mu_*^+ (\|\hat{x}_*\|^2 - \delta^2) = 0. \quad (18)$$

Now compare \bar{x}_* and μ_* from (10)–(12) with \hat{x}_* and μ_*^+ from (15), (17) and (18). The only substantial difference is between (11) and (15). Indeed, if $\mu_*^+ \geq \max(0, -\lambda_1)$, the two sets of conditions are identical, and in this case $x_m = x_*$ and $\mu_m = \mu_*$, i.e. the solution from the subspace \mathcal{K}_m solves the full-space trust-region problem (1). This must occur if $\min_{j \in \mathcal{I}_+} \lambda_j = \lambda_1$ or, equivalently $\mathcal{I}_+ \cap \{1, \dots, n_1\} \neq \emptyset$, where we recall n_1 is the multiplicity of λ_1 , but may also happen if $\min_{j \in \mathcal{I}_+} \lambda_j > \lambda_1$. If $\mu_*^+ < -\lambda_1$, $\mathcal{I}_+ \cap \{1, \dots, n_1\} = \emptyset$, and x_m cannot solve (1), but it is nonetheless a critical point of the problem.⁴ This possibility is often called the ‘hard case’ [30] and $\mu_* = -\lambda_1$; the first block equation in (10) uniquely defines \bar{x}_*^+ , and \bar{x}_*^0 is a multiple of any eigenvector of the second (singular) block, the precise combination ensuring that $\|\bar{x}_*\| = \delta$.

The main lesson here is that if we wish to solve (1) we shall have to look outside the Krylov space and may need to compute an eigenvector corresponding to λ_1 . If we are content simply in finding a critical point of (1), the Krylov space suffices. An essentially identical argument may be used in the case of the regularization subproblem (2) with the same conclusions.

3. Error bounds

3.1. Bounds on the decrease of the objective functions

In essence Carmon and Duchi [8] provide the following bounds.⁵

Theorem 3.1 ([8, Theorem 1 & Corollary 3]): *Let λ_1 and λ_n be the leftmost and rightmost eigenvalues of H . Then, for all $k \geq 0$,*

(i)

$$q(x_k) - q(x_*) \leq 36[q(0) - q(x_*)] \left(e^{-4\sqrt{(\lambda_1 + \mu_*)/(\lambda_n + \mu_*)}} \right)^k, \quad (19)$$

where x_k is given by (4), x_* is a global minimizer of (1), and μ_* is its corresponding Lagrange multiplier, and

(ii)

$$q^R(x_k, \sigma, p) - q^R(x_*, \sigma, p) \leq 36 [q^R(0, \sigma, p) - q^R(x_*, \sigma, p)] \left(e^{-4\sqrt{(\lambda_1 + \mu_*)/(\lambda_n + \mu_*)}} \right)^k, \quad (20)$$

where x_k is given by (5), x_* is a global minimizer of (2) and $\mu_* = \sigma \|x_*\|$.

Thus the error in the relevant objective function decreases at worst linearly as a function of the subspace dimension unless $\mu_* = -\lambda_1$, in which case Theorem 3.1 provides no useful bound. As we have already mentioned, the unlikely ‘hard case’ $\mu_* = -\lambda_1$ only occurs if g is orthogonal to the space of eigenvectors corresponding to the eigenvalue λ_1 of H , and should this happen these eigenvectors will not occur in the Krylov spaces \mathcal{K}_k , except through numerical rounding. A simple expedient advocated by others [8] is to perturb g by a small random vector so that, with high probability, it will have some component in the space of eigenvectors corresponding to the eigenvalue λ_1 , and thus the hard case cannot occur.

We note that Carmon and Duchi actually provide a second, sublinear decrease estimate that may be less pessimistic for small k , but we shall not use this here. Zhang et al. [42, Theorem 4.2] suggest quantitatively-similar bounds in the trust-region case, but again we shall not need them.

We now restrict our attention to the best estimate x_m available from the evolving Krylov space. We provide a bound on the decrease of the objective function at this point compared to that at the best estimate x_k available from \mathcal{K}_k ; the bound indicates at worst a linear rate of convergence as the space expands. We exclude the special case $g = 0$ since then $x = 0$ is a critical point of both of the subproblems under consideration.

Corollary 3.2: *Suppose that $g \neq 0$. Let $\{\lambda_j, u_j\}_{j \in \mathcal{N}}$ be eigenpairs of H , $\mathcal{I}_+ = \{j \mid g^T u_j \neq 0\}$, $m = |\mathcal{I}_+|$, and*

$$\lambda_+^{\min} = \min_{j \in \mathcal{I}_+} \lambda_j \quad \text{and} \quad \lambda_+^{\max} = \max_{j \in \mathcal{I}_+} \lambda_j. \quad (21)$$

Then, for all $k \geq 0$,

(i)

$$q(x_k) - q(x_m) \leq 36 [q(0) - q(x_m)] \left(e^{-4/\sqrt{\kappa_m}} \right)^k, \quad (22)$$

where x_k and x_m are given by (4),

$$\kappa_m := \frac{\lambda_+^{\max} + \mu_m}{\lambda_+^{\min} + \mu_m}, \quad (23)$$

and μ_m is the Lagrange multiplier corresponding to x_m , and

(ii)

$$q^R(x_k, \sigma, p) - q^R(x_m, \sigma, p) \leq 36 [q^R(0, \sigma, p) - q^R(x_m, \sigma, p)] \left(e^{-4/\sqrt{\kappa_m}} \right)^k, \quad (24)$$

where x_k and x_m are given by (5), κ_m is given by (23) but now $\mu_m = \sigma \|x_m\|$.

Proof: Since $g \neq 0$, $\mathcal{I}_+ \neq \emptyset$, $m > 0$ and both λ_+^{\min} and λ_+^{\max} are well defined. Let H_+ be the matrix with eigenpairs $\{\lambda_j, u_j\}$ for $j \in \mathcal{I}_+$ and $\{\lambda_+^{\max}, u_j\}$ for $j \in \mathcal{I}_0 = \mathcal{N} \setminus \mathcal{I}_+$. Then $\mathcal{K}_k = \text{span}\{H^i g\}_{i=0}^{k-1} = \text{span}\{H_+^i g\}_{i=0}^{k-1}$ for $k \geq 0$, and the iterates x_k generated from the Krylov spaces \mathcal{K}_k for (4) and (5) for the problem with Hessian H_+ are identical to those with Hessian H . However the hard case cannot occur with the Hessian H_+ as none of the eigenvalues for $j \in \mathcal{I}_0$ is smaller than the smallest for $j \in \mathcal{I}_+$. Hence, as we saw in §2, for this Hessian $x_* = x_m$ and $\mu_* = \mu_m$. Thus we may apply Theorem 3.1 for the problem with Hessian H_+ to deduce (22) and (24). ■

As Carmon and Duchi mention, this then implies a worst-case estimate of

$$k \leq \min(m, O(\sqrt{\kappa_m} \log(1/\epsilon))) \quad (25)$$

iterations in order to guarantee $q(x_k) - q(x_m) \leq \epsilon$ or $q^R(x_k, \sigma, p) - q^R(x_m, \sigma, p) \leq \epsilon$ as appropriate.

3.2. Bounds on the residuals

Recall that the orthonormal Lanczos basis matrix $V_k \in \mathfrak{R}^{n \times k}$ for \mathcal{K}_k satisfies

$$HV_k = V_k T_k + \gamma_k v_{k+1} e_k^T, \quad (26)$$

where

$$T_k = \begin{pmatrix} \delta_1 & \gamma_1 & & & \\ \gamma_1 & \delta_2 & & & \\ & & \cdot & & \\ & & & \cdot & \\ & & & & \delta_{k-1} & \gamma_{k-1} \\ & & & & \gamma_{k-1} & \delta_k \end{pmatrix} \quad (27)$$

is tridiagonal and the γ_i , $i = 1, \dots, k-1$ with $k \leq m$, are strictly positive [2,22,27,39]. As the off diagonals $\gamma_i > 0$, T_k is irreducible, and has distinct real eigenvalues (the so-called Ritz values) $\theta_{i,k}$, $i = 1, \dots, k$, arranged in increasing order. It is well known [17,

Corollary 8.1.7] that the Ritz values satisfy the interlacing properties

$$\theta_{i,k} \leq \theta_{i,k+1} \leq \theta_{i+1,k} \quad (28)$$

for $i = 1, \dots, k, k < m$, and

$$\lambda_+^{\min} = \theta_{1,m} \leq \theta_{1,k+1} \leq \theta_{1,k} \leq \theta_{1,1} \equiv \frac{g^T H g}{\|g\|^2} \leq \theta_{k,k} \leq \theta_{k+1,k+1} \leq \theta_{m,m} = \lambda_+^{\max} \quad (29)$$

for $k = 1, \dots, m - 1$, where λ_+^{\min} and λ_+^{\max} are defined in (21).

Since $v_1 = g/\|g\|$, it follows that

$$V_k^T g = \|g\| e_1 \quad \text{and} \quad g = \|g\| V_k e_1 \quad (30)$$

as V_k has orthonormal columns. Furthermore pre-multiplying (26) by V_k^T yields

$$V_k^T H V_k = T_k,$$

and thus the definition (4) implies that

$$x_k = V_k y_k, \quad \text{where} \quad (T_k + \mu_k I) y_k = V_k^T (H + \mu_k I) V_k y_k = -V_k^T g = -\|g\| e_1. \quad (31)$$

Moreover, applying Theorem 2.1 to (4) shows that $T_k + \mu_k I$ is positive definite.

Let

$$r_k := g + H x_k + \mu_k x_k = g + (H + \mu_k I) x_k. \quad (32)$$

It then follows from (26), (30) and (31) that

$$\begin{aligned} H x_k &= H V_k y_k = V_k T_k y_k + \gamma_k e_k^T y_k v_{k+1} = -V_k (\mu_k y_k + \|g\| e_1) + \gamma_k e_k^T y_k v_{k+1} \\ &= -\mu_k x_k - g + \gamma_k e_k^T y_k v_{k+1}. \end{aligned}$$

Hence $r_k = \gamma_k e_k^T y_k v_{k+1}$ and

$$\|r_k\| = \gamma_k \left| e_k^T y_k \right| = \gamma_k \|g\| \left| e_k^T (T_k + \mu_k I)^{-1} e_1 \right|. \quad (33)$$

Note that the definition of $\gamma_k > 0$ as the $(k, k + 1)$ -st entry of T_{k+1} and the Cauchy Schwarz inequality implies that

$$\gamma_k = e_{k+1}^T T_{k+1} e_k \leq \|T_{k+1}\| = \|V_{k+1}^T H V_{k+1}\| \leq \|H\|. \quad (34)$$

Thus, aside from the term $\gamma_k \|g\| > 0$, the residual norm decays with $|e_k^T (T_k + \mu_k I)^{-1} e_1|$, and we now focus on this.

We recall a vital result by Demko et al. [12] on the component-wise decay of the inverse of symmetric banded matrices. Here the bandwidth of a banded symmetric matrix M is the number of nonzero upper (or equivalently lower) super diagonals, and, if M is additionally positive definite, $\kappa(M) := \lambda_{\max}(M)/\lambda_{\min}(M)$ is its spectral condition number, where $0 < \lambda_{\min}(M) \leq \lambda_{\max}(M)$ are the left- and right-most eigenvalues of M .

Lemma 3.3 ([12, Theorem 2.4]): *Let $M \in \mathbb{R}^{n \times n}$ be a symmetric, positive definite matrix with bandwidth $\beta > 0$. Then*

$$|(M^{-1})_{ij}| \leq ct^{i-j/\beta}$$

for all $i, j = 1, \dots, n$, where

$$c = \frac{1}{\lambda_{\min}(M)} \max \left(1, \frac{(\sqrt{\kappa(M)} + 1)^2}{2\kappa(M)} \right) \quad \text{and} \quad t = \frac{\sqrt{\kappa(M)} - 1}{\sqrt{\kappa(M)} + 1}.$$

Note that we shall prefer the slightly weaker, but simpler, bound

$$c \leq \frac{2}{\lambda_{\min}(M)}, \quad (35)$$

and indeed $c = 1/\lambda_{\min}(M)$ so long as $\kappa(M) \geq \sqrt{1 + \sqrt{2}}$.

For any $k \leq m$, we have that T_k from (27) is symmetric and tridiagonal with left- and right-most eigenvalues (Ritz values) respectively $\theta_{1,k} < \theta_{k,k}$. As we have seen $T_k + \mu_k I$ is positive definite, and thus has distinct left- and right-most eigenvalues

$$\lambda_{\min}(T_k + \mu_k I) \equiv \theta_{1,k} + \mu_k < \lambda_{\max}(T_k + \mu_k I) \equiv \theta_{k,k} + \mu_k \quad (36)$$

as well as spectral condition number

$$\kappa_k := \kappa(T_k + \mu_k I) = \frac{\theta_{k,k} + \mu_k}{\theta_{1,k} + \mu_k}. \quad (37)$$

We may apply Lemma 3.3 to $T_k + \mu_k I$ to deduce our main result.

Theorem 3.4: *The residual (32) for the k th iterate, x_k^* , generated by either the trust-region subproblem (4) or the regularization subproblem (5) satisfies the bound*

$$\|r_k\| \leq \|g\| \left(\frac{2\gamma_k \kappa_k}{\theta_{k,k} + \mu_k} \right) \left(\frac{\sqrt{\kappa_k} - 1}{\sqrt{\kappa_k} + 1} \right)^{k-1}, \quad (38)$$

where κ_k is given by (37) and γ_k is the $(k, k+1)$ -st entry of T_{k+1} .

Proof: Since $|e_k^T (T_k + \mu_k I)^{-1} e_1| = |((T_k + \mu_k I)^{-1})_{k,1}|$, we may apply Lemma 3.3 to $T_k + \mu_k I$, with $\beta = 1$, together with (35)–(37) to deduce the bound

$$\left| e_k^T (T_k + \mu_k I)^{-1} e_1 \right| \leq c_k t_k^{k-1} \quad (39)$$

for all $k \leq m$, where

$$c_k = \frac{2}{\theta_{1,k} + \mu_k} \equiv \frac{2\kappa_k}{\theta_{k,k} + \mu_k} \quad \text{and} \quad t_k = \frac{\sqrt{\kappa_k} - 1}{\sqrt{\kappa_k} + 1}. \quad (40)$$

The desired bound (38) then follows directly from (33) and (40). ■

In the trust-region case, this leads to the following residual bound.⁶

Corollary 3.5: *The residual (32) for the k th iterate, x_k^* , generated by the trust-region subproblem (4) satisfies the bound*

$$\|r_k\| \leq \|g\| \left(\frac{2\delta \|H\| \kappa_k}{\|g\|} \right) \left(\frac{\sqrt{\kappa_k} - 1}{\sqrt{\kappa_k} + 1} \right)^{k-1}, \quad (41)$$

where κ_k is given by (37).

Proof: It follows from (31), the Cauchy-Schwarz and Rayleigh-Ritz inequalities and the bound $\|y_k\| \leq \delta$ that

$$\|g\|^2 = \|(T_k + \mu_k I)y_k\|^2 \leq (\theta_{k,k} + \mu_k)^2 \|y_k\|^2 \leq (\theta_{k,k} + \mu_k)^2 \delta^2,$$

and hence

$$\frac{1}{\theta_{k,k} + \mu_k} \leq \frac{\delta}{\|g\|}. \quad (42)$$

Thus combining (34), (38) and (42), we find that (41) holds. ■

Corollary 3.6: *Let $\kappa_* = \max_{1 \leq k \leq m} \kappa_k$. Then the iteration defined by (4) satisfies*

$$\|r_k\| \leq \epsilon \quad (43)$$

as soon as

$$k \leq \min \left[m, \left\lceil \log \left(\frac{2\delta \|H\| \kappa_*}{\epsilon} \right) / \log \left(\frac{\sqrt{\kappa_*} + 1}{\sqrt{\kappa_*} - 1} \right) \right\rceil + 1 \right]. \quad (44)$$

Proof: Since the function

$$q(\kappa) = \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1}$$

is monotonically increasing for $\kappa \geq 1$, we deduce from Corollary 3.5 that

$$\|r_k\| \leq \|g\| \left(\frac{2\delta \|H\| \kappa_*}{\|g\|} \right) \left(\frac{\sqrt{\kappa_*} - 1}{\sqrt{\kappa_*} + 1} \right)^{k-1}. \quad (45)$$

Recalling that $r_m = 0$, (43) and (45) lead directly to (44). ■

A similar result is possible for the regularization case.

Corollary 3.7: *The residual (32) for the k th iterate, x_k^* , generated by the regularization subproblem (5) satisfies the bound*

$$\|r_k\| \leq \|g\| \left(\frac{2\|H\| \kappa_k}{\|g\|} \right) \left(\frac{\mu_k}{\sigma} \right)^{1/(p-2)} \left(\frac{\sqrt{\kappa_k} - 1}{\sqrt{\kappa_k} + 1} \right)^{k-1}, \quad (46)$$

where κ_k is given by (37).

Proof: It follows from (31), the Cauchy-Schwarz and Rayleigh-Ritz inequalities and the relationship $\mu_k = \sigma \|y_k\|^{p-2}$ that

$$\|g\|^2 = \|(T_k + \mu_k I)y_k\|^2 \leq (\theta_{k,k} + \mu_k)^2 \|y_k\|^2 = (\theta_{k,k} + \mu_k)^2 \left(\frac{\mu_k}{\sigma}\right)^{2/(p-2)}$$

and hence

$$\frac{1}{\theta_{k,k} + \mu_k} \leq \frac{1}{\|g\|} \left(\frac{\mu_k}{\sigma}\right)^{1/(p-2)}. \quad (47)$$

Thus (46) follows by combining (34), (38) and (47). \blacksquare

Corollary 3.8: Let $\kappa_* = \max_{1 \leq k \leq m} \kappa_k$. Then the iteration defined by (4) satisfies (43) as soon as

$$k \leq \min \left[m, \left\lceil \log \left(\frac{2\|H\|\kappa_*}{\epsilon} \left(\frac{\mu_m}{\sigma}\right)^{1/(p-2)} \right) / \log \left(\frac{\sqrt{\kappa_*} + 1}{\sqrt{\kappa_*} - 1} \right) \right\rceil + 1 \right]. \quad (48)$$

Proof: Given Corollary 3.7, the proof is essentially identical to that of Corollary 3.6, except that we additionally make use of the bound $0 \leq \mu_k \leq \mu_m$ for $k = 1, \dots, m$ [10, Theorem 2.5]. \blacksquare

3.3. Comments

We now comment on the bounds obtained in §3.1 and 3.2. Those on the (linear) rates of convergence given in Corollaries 3.2, 3.6 and 3.8 are very typical of the Chebyshev bounds that have been derived for conjugate-gradient (CG)-like methods for solving symmetric, positive-definite systems of linear equations $Ax = b$ (see, e.g. [27, §5.6.2]). Briefly, in this case $\|r_k\|_{A^{-1}} = \|x_k - x_*\|_A$, where $r_k = Ax_k - b$ and $x_* = A^{-1}b$. Since $\|r_k\| \leq \sqrt{\lambda_{\max}(A)} \|r_k\|_{A^{-1}}$, the argument in the CG case focuses on finding an upper bound on $\|x_k - x_*\|_A$. In particular, x_k is chosen to minimize $\|x - x_*\|_A$ over all $x \in \mathcal{K}(A, b, k)$, and this is easily shown to lead to the bound

$$\|x_k - x_*\|_A \leq \|x_0 - x_*\|_A \min_{\psi \in \mathcal{P}_k} \max_{i \in \mathcal{N}} |\psi(\lambda_i(A))|, \quad (49)$$

where

$$\mathcal{P}_k = \{\text{polynomials } \psi \text{ of degree } k \text{ for which } \psi(0) = 1\}$$

and $\lambda_i(A)$, $i \in \mathcal{N}$, are a subset of the eigenvalues of A . Weakening the requirement that the maximum in (49) instead considers $\psi(\lambda)$ over all $\lambda \in [\lambda_{\min}(A), \lambda_{\max}(A)]$ and invoking a well-known bound from approximation theory relating to Chebyshev polynomials, it then follows that

$$\|x_k - x_*\|_A \leq 2 \left(\frac{\sqrt{\kappa(A)} - 1}{\sqrt{\kappa(A)} + 1} \right)^k \|x_0 - x_*\|_A.$$

Since $\|x_0 - x_*\|_A \leq \|r_0\| / \sqrt{\lambda_{\min}(A)}$, we thus obtain the bound

$$\|r_k\| \leq 2\sqrt{\kappa(A)} \left(\frac{\sqrt{\kappa(A)} - 1}{\sqrt{\kappa(A)} + 1} \right)^k \|r_0\|;$$

if $x_0 = 0$, $\|r_0\| = \|b\|$ in the latter.

The presence of $\kappa_m \leq \kappa_*$ in the bounds in §3.1 and 3.2 is strongly reminiscent of the CG case, and indicates that rescaling (preconditioning) the problem so that κ_m or $\kappa_* = O(1)$ would be beneficial. In the strictly convex case when H is positive definite, κ_m is no larger than the traditional condition number $\lambda_+^{\max}/\lambda_+^{\min}$ obtained from (29). Although we know that $\theta_{k,k}$ increases monotonically from (29), as does μ_k [10,28], we have not been able to prove that κ_k increases monotonically,⁷ albeit in practice it appears to.

We need to be very cautious here as although such bounds accurately predict the worst possible case [8,22], they are often very pessimistic in general, a point stressed in [27]. We shall return to this in §4. Nevertheless, if one is interested in the worst-case, bounds such as (25), (44) and (48) are relevant.

We tried two other approaches to derive useful bounds on the norm of the residual, (32). The first aims to use the known decrease in the model objective given by Corollary 3.2 to deduce a similar bound on $\|r_k\|$. Since x_m from (4) is a critical point of (1), we have that

$$g + Hx_m + \mu_m x_m = 0, \quad (50)$$

and hence

$$\begin{aligned} r_k &:= g + Hx_k + \mu_k x_k = -Hx_m - \mu_m x_m + Hx_k + \mu_k x_k \\ &= (H + \mu_m I)(x_k - x_m) + (\mu_k - \mu_m)x_k. \end{aligned} \quad (51)$$

Elementary manipulation of these (see Appendix 1) then leads to the bound

$$\begin{aligned} (\lambda_+^{\max} + \mu_m)^{-1} \|r_k\|^2 &\leq 2 [q(x_k) - q(x_m)] + \rho_k, \quad \text{where} \\ \rho_k &:= \mu_k (\|x_k\|^2 - \|x_m\|^2) - (\mu_m - \mu_k) \|x_m - x_k\|^2 + (\mu_m - \mu_k)^2 x_k^T (H + \mu_m I)^{-1} x_k, \end{aligned} \quad (52)$$

which exposes the dependence on the model objective decrease $q(x_k) - q(x_m)$. Unfortunately, aside from the case where $\mu_m = 0$ for which $\rho_k = 0$, we are not able to find a useful bound on ρ_k ; ideally we would like to show that $\rho_k \leq 0$. Of course, even had we had succeeded in bounding ρ_k , the overall bound we would have obtained via Corollary 3.2 for the $q(x_k) - q(x_m)$ term would not have been substantially different from those given in Corollaries 3.6 and 3.8.

Our second approach tried to mimic that taken for the CG method for positive-definite linear systems. However, the argument relating the A -norm of the error to the A^{-1} -norm of the residual and the subsequent min-max characterization depends crucially on the definiteness of A , and thus this line of attack is not obvious for our case where H may be indefinite. Nevertheless it is easy to derive the bound

$$\|r_k\| \leq \max_{j \in \mathcal{I}_+} |\psi_k(\lambda_j + \mu_k)| \|g\|, \quad \text{where} \quad \psi_k(\lambda_j + \mu_k) = \prod_{i=1}^k \frac{(\theta_{i,k} - \lambda_j)}{(\theta_{i,k} + \mu_k)}. \quad (53)$$

on the residual (32) (see Appendix 2). Although we do not know how to derive a useful bound on $\psi_k(\lambda_j + \mu_k)$, as we see in §4, to do so might provide a much closer match to the true residual than provided by Corollaries 3.6 and 3.8.

4. Experiments

We consider nine examples that aim to illustrate our analysis; all nine are available as part of the CUTEst [21] set of test problems. Each is of the form

$$q(x) = \frac{1}{2} \sum_{i=1}^n d_i x_i^2 + \sum_{i=1}^n x_i,$$

but vary according to the diagonal Hessian elements, d_i . Specifically we have examples

$$\begin{aligned} \text{DIAGPQT: } & d_i = -i^2/n + n + 1/n, \\ \text{DIAGPQE: } & d_i = i, \\ \text{DIAGPQB: } & d_i = i^2/n, \\ \text{DIAGIQT: } & d_i = -i^2/n + n/2 + 1/n, \\ \text{DIAGIQE: } & d_i = i - n/2, \\ \text{DIAGIQB: } & d_i = i^2/n - n/2 + 1/n, \\ \text{DIAGNQT: } & d_i = -i^2/n, \\ \text{DIAGNQE: } & d_i = i - n - 1, \text{ and} \\ \text{DIAGNQB: } & d_i = i^2/n - n - 1/n, \end{aligned}$$

for $i = 1, \dots, n$; in our tests we let $n = 1000$, and ignore the additional simple-bound constraints specified in the CUTEst examples. The first three are convex with Hessian spectra that bunch towards the bottom of the range, that are equispaced, and that coalesce towards the top of the range respectively. The second three shift the spectra of the first three downwards by roughly $n/2$, leading to indefinite Hessians, while the last three concave examples shift downwards by more or less $n + 1/n$.

In Figure 1 we compare the true residual against bounds derived in Section 3 when running GALAHAD's [19] GLTR package [18] to solve the trust-region subproblem (1) on the first three test examples. Almost identical plots have been obtained for the remaining examples for the trust-region case since the residual

$$r_k = Hx_k + \mu_k x_k + g = (H - \omega I)x_k + (\mu_k + \omega)x_k + g$$

shows that shifting the Hessian downwards by ω is compensated by shifting the multiplier upwards by the same amount once the trust-region constraint is active.

We observe that although Theorem 3.4 and Corollary 3.5 provide bounds on the residual, they may be far from sharp, especially when the spectrum is equispaced or bunched towards the top end. In particular, the bounds do not capture the superlinear behaviour of the residuals in these cases; the slopes best mimic those from the earlier iterations. This largely agrees with the observations made and conclusions drawn in the linear-equation case [27]. The inferiority of the bound in Corollary 3.5 compared to that in Theorem 3.4 merely reflects the weakening that results when approximating unknown quantities (i.e. $\theta_{k,k} + \mu_k$) by known ones (i.e. H, g, δ). By contrast, the bound provided by (53) is quantitatively far better, but, of course, this requires full knowledge of the spectrum.

Figures 2–4 compare the estimates (38), (46) and (53) against the true residual when running GALAHAD's GLRT package [9] to solve the regularization subproblem (2) on all nine test examples; unlike for the trust-region case, a translation of the Hessian does not

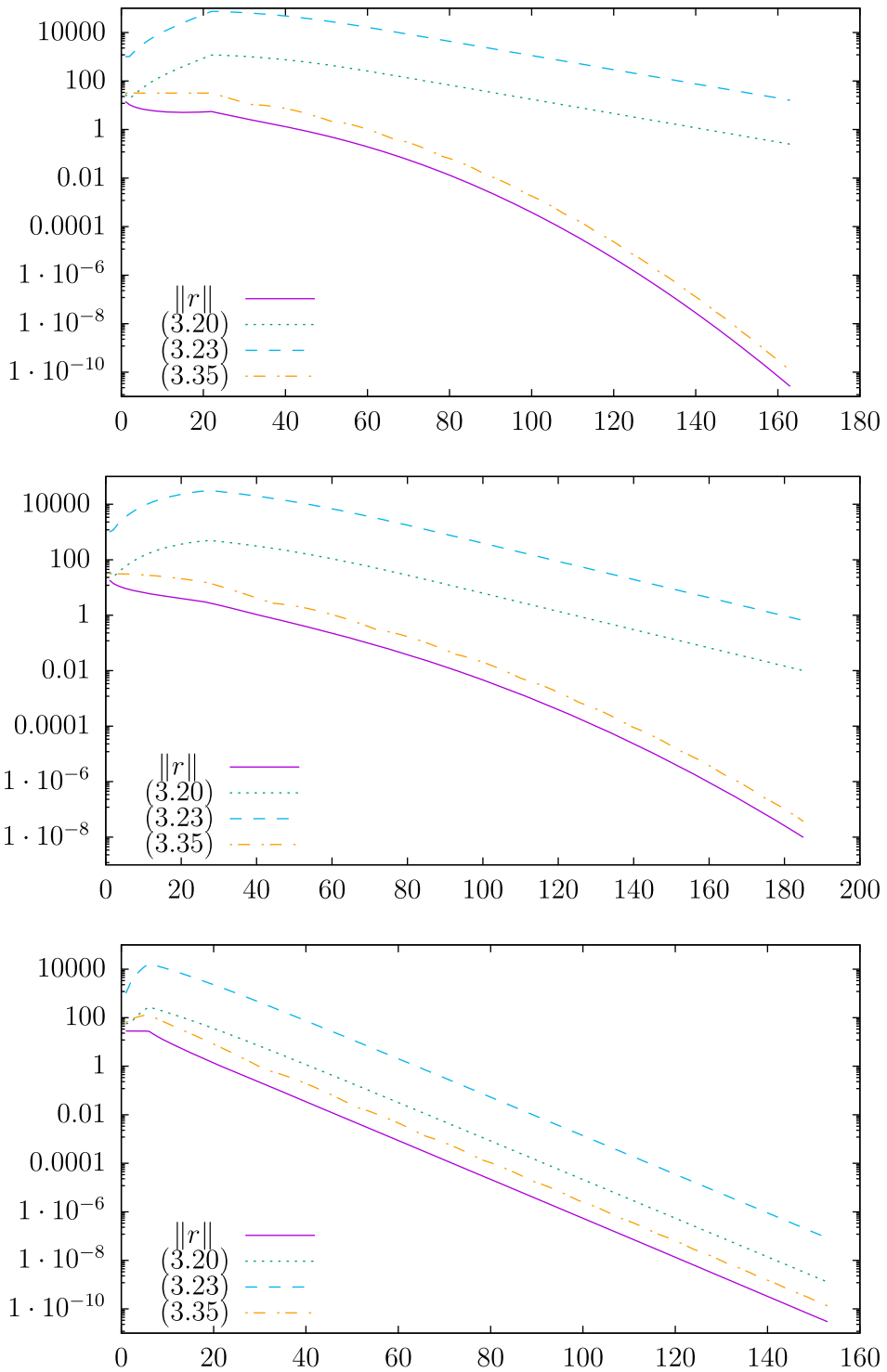


Figure 1. \log_{10} of the residual (y-axis) as the iteration proceeds (x-axis) for GLTR as applied to the convex problems DIAGPQT (top plot), DIAGPQE (middle) and DIAGPQB (bottom) with a trust-region radius $\delta = 1$. Each figure shows the residual (33) (solid line), and the estimates (38) (dotted line), (41) (dashed line) and (53) (dash-dot line).

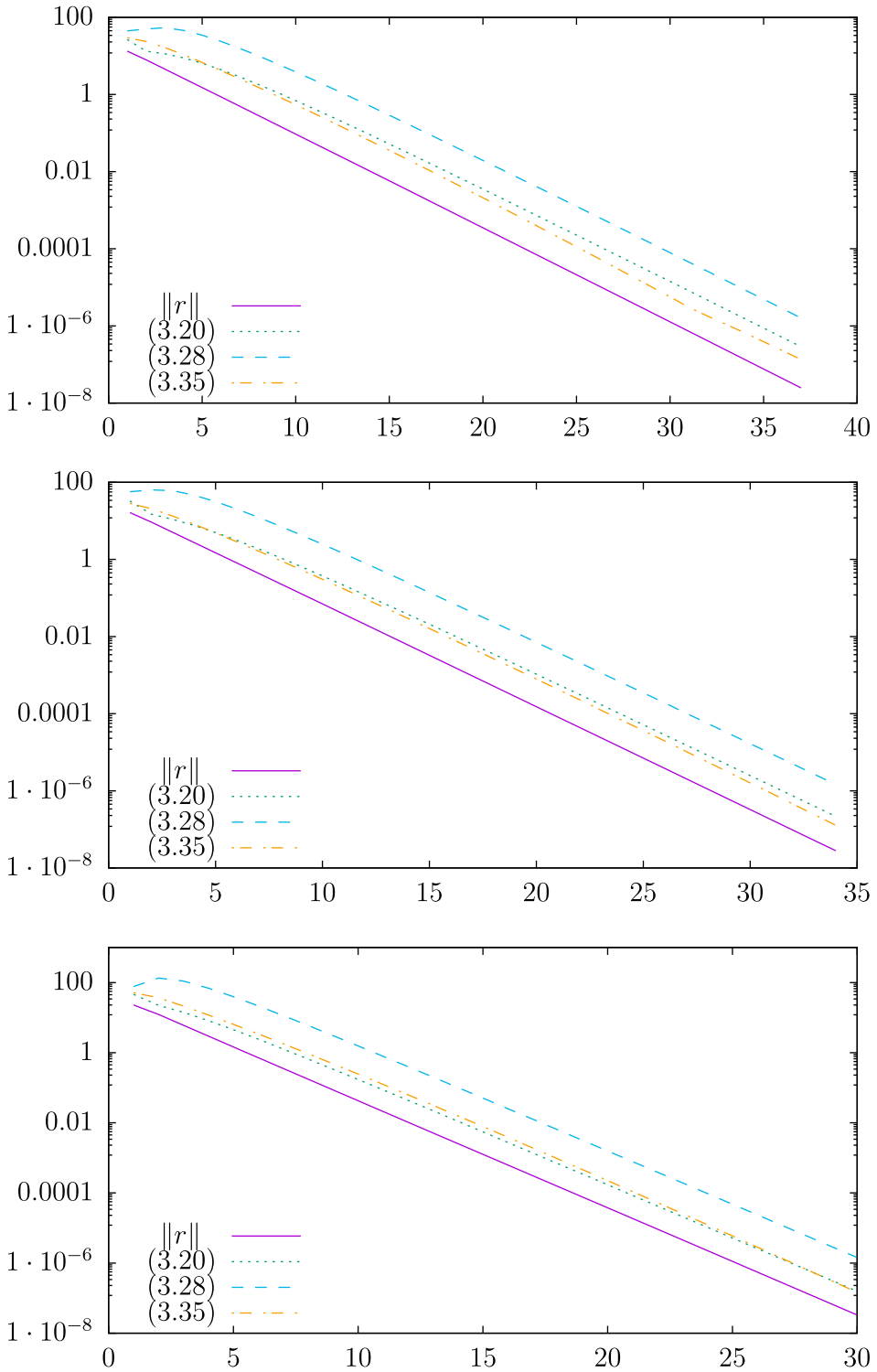


Figure 2. \log_{10} of the residual (y-axis) as the iteration proceeds (x-axis) for GLRT as applied to the convex problems `DIAGPQT` (top plot), `DIAGPQE` (middle) and `DIAGPQB` (bottom) with a regularization weight $\sigma = 1000$ and $p = 3$. Each figure shows the residual (33) (solid line), and the estimates (38) (dotted line), (46) (dashed line) and (53) (dash-dot line).

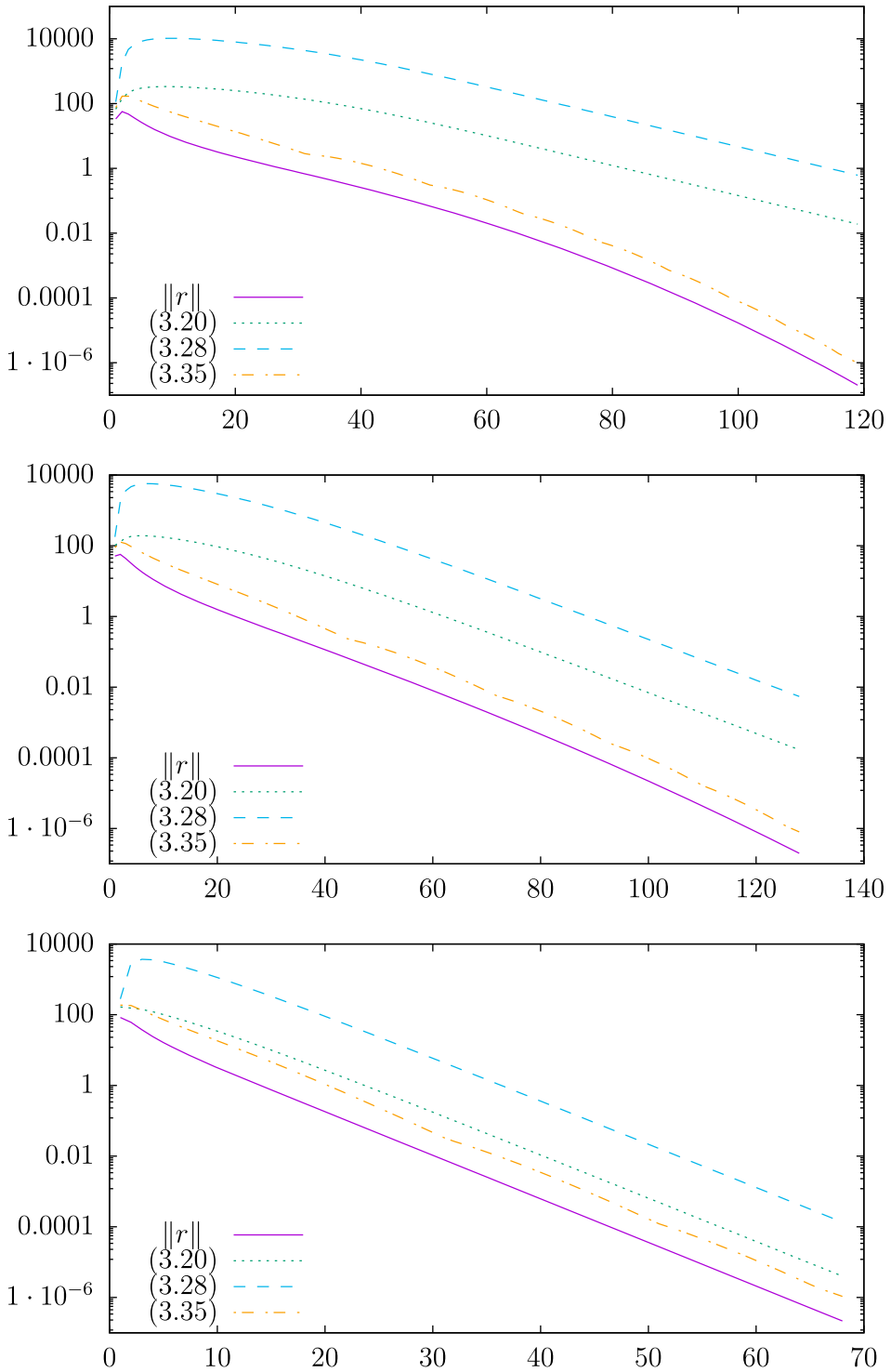


Figure 3. \log_{10} of the residual (y-axis) as the iteration proceeds (x-axis) for GLRT as applied to the indefinite problems `DIAGIQT` (top plot), `DIAGIQE` (middle) and `DIAGIQB` (bottom) with a regularization weight $\sigma = 1000$ and $p = 3$. Each figure shows the residual (33) (solid line), and the estimates (38) (dotted line), (46) (dashed line) and (53) (dash-dot line).

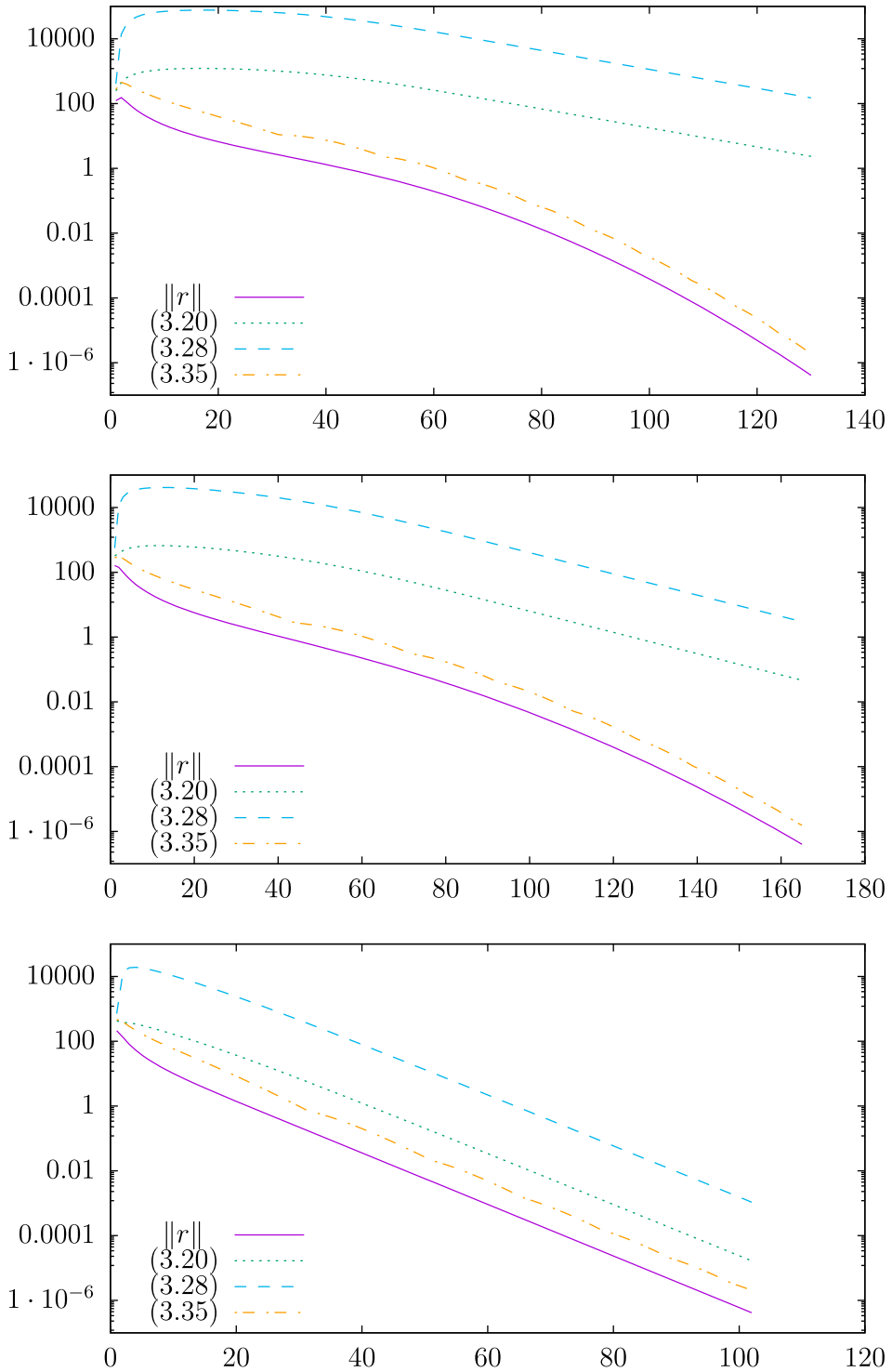


Figure 4. \log_{10} of the residual (y-axis) as the iteration proceeds (x-axis) for GLRT as applied to the concave problems DIAGNQT (top plot), DIAGNQE (middle) and DIAGNQB (bottom) with a regularization weight $\sigma = 1000$ and $p = 3$. Each figure shows the residual (33) (solid line), and the estimates (38) (dotted line), (46) (dashed line) and (53) (dash-dot line).

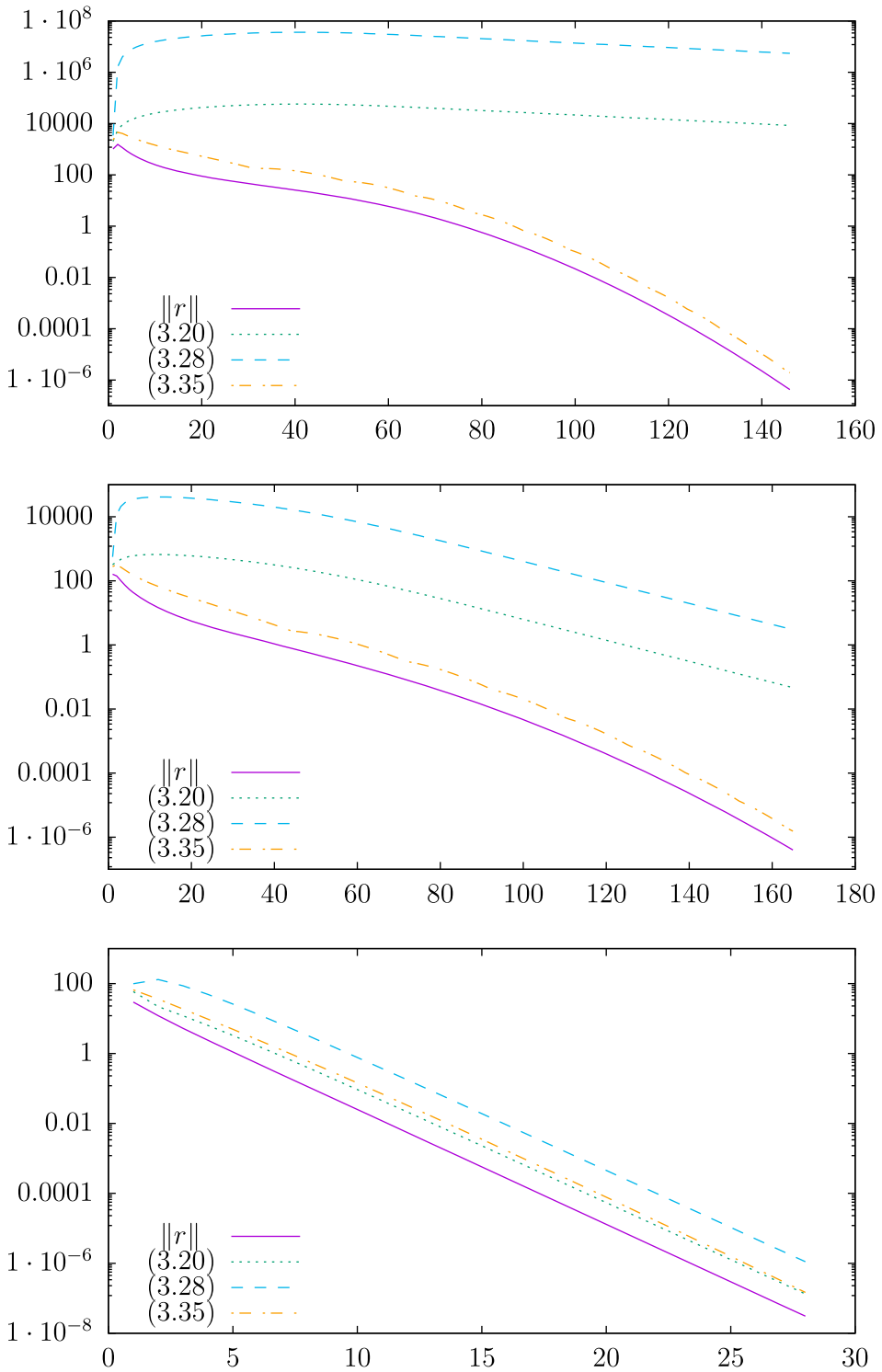


Figure 5. \log_{10} of the residual (y-axis) as the iteration proceeds (x-axis) for GLRT as applied to the concave problem `DIAGNQT` with different values of σ when $p = 3$. Specifically $\sigma = 100$ (top plot), $\sigma = 1000$ (middle) and $\sigma = 10000$ (bottom). Each figure shows the residual (33) (solid line), and the estimates (38) (dotted line), (46) (dashed line) and (53) (dash-dot line).

produce essentially identical plots when moving from the convex via the indefinite to the concave cases.

Once again, we observe that the bounds (38), (46) may be far from sharp, and can fail to reflect the later superlinear convergence of the residuals. The behaviour is most extreme for the concave examples, and for those whose spectra coalesce at the top of their ranges. As before, (53) provides a much more faithful bound.

Finally Figures 5 illustrate the effect of changing the regularization weight, σ , when solving (2) for the example DIAGNQT, on the residual estimates. The subproblems become increasingly hard as σ shrinks, and the estimates correspondingly poorer. Indeed, the decrease predicted by (38) when $\sigma = 100$ barely indicates convergence, while although the rates for the actual residual and the prediction (53) are initially slow, they later accelerate.

5. Conclusions

We have derived bounds for the objective errors and gradient residuals when finding approximations to the solution of common regularized quadratic optimization problems within evolving Krylov spaces. Those for the objective errors are trivial extensions of existing ones [8], while those for the gradient residuals generalize well-known ones for conjugate-gradient methods applied to definite linear systems. Quantitatively the bounds behave just as in the conjugate-gradient case, but reflect additional complication of the subproblems, particularly the potential indefiniteness of the matrices involved.

We express some caution since in exceptional cases Krylov methods may not find the global solutions to our problems. If this is the goal, then additional precautions [8,18] that are not covered by our bounds may be necessary. If our goal is simply to find an approximation that yields a small gradient residual—and this is often the case when the subproblem occurs as component of a more general optimization calculation—then our bounds are appropriate, and provide upper bounds on the number of iterations required to achieve a given stated accuracy.

Our bounds do not reflect the ‘superlinear’ behaviour that is sometimes observed in practice that results from annihilation of extreme eigenvalues by the Krylov process. A more sophisticated analysis, akin to that by Axelsson, Kaporin and others [3,4], might provide this, but we have not attempted it.

Notes

1. That the ‘arg min’s in (4) and (5) are unique follows from Theorems 2.1 and 2.2, since g lies in \mathcal{K}_k by construction.
2. The precise details of the implementations of GLTR and GLRT are, of course, important, but of no consequence in the bounds that we derive. Such bounds apply equally for any method that use the iterates (4) or (5).
3. This is equivalently the *grade* of H with respect to g .
4. It will only be a local minimizer of $\mu_*^+ > -\lambda_2$ [29].
5. Strictly [8, Corollary 3] only considers the case $p = 3$, but their method of proof holds in general.
6. Another thing we know but we have not used.: $0 \leq \mu_1 \leq \mu_k \leq \mu_m$ for $k = 1, \dots, m$ [28].
7. The result would follow if we could show that $\theta_{1,k} + \mu_k$ decreases monotonically with growing k .

Acknowledgments

The authors are very grateful for fruitful discussions on aspects of this work with Coralia Cartis, Tyrone Rees and Philippe Toint. We also appreciate very helpful comments from two referees and an associate editor.

Disclosure statement

No potential conflict of interest was reported by the authors.

Funding

This work was supported by the Engineering and Physical Sciences Research Council (EPSRC) grant EP/M025179/1.

Notes on contributors

Nicholas I. M. Gould, an STFC Senior Fellow and Professor of Numerical Optimization, has interests in theoretical and practical optimization, and links to linear algebra. He is a SIAM Fellow, a winner of the Leslie Fox and Beale-Orchard Hays prizes, and his most recent book is *Trust-region Methods* (SIAM, 2000), written jointly with Andy Conn and Philippe Toint.

Valeria Simoncini, a Professor of Numerical Analysis, has interests in matrix analysis and computations with applications to scientific computing and engineering. She is a SIAM Fellow, and member of several editorial board, including SIAM J. Matrix Analysis and Applications, and SIAM J. Numerical Analysis.

ORCID

Nicholas I. M. Gould  <http://orcid.org/0000-0002-1031-1588>

Valeria Simoncini  <http://orcid.org/0000-0003-0795-5865>

References

- [1] S. Adachi, S. Iwata, Y. Nakatsukasa, and A. Takeda, *Solving the trust-region subproblem by a generalized eigenvalue problem*, SIAM J. Optim. 27(1) (2017), pp. 269–291.
- [2] O. Axelsson, *Iterative Solution Methods*, Cambridge University Press, Cambridge, 1996.
- [3] O. Axelsson and I. Kaporin, *On the sublinear and superlinear rate of convergence of conjugate gradient methods*, Numer. Algorithms 25(1) (2000), pp. 1–22.
- [4] O. Axelsson and J. Karátson, *Reaching the superlinear convergence phase of the CG method*, J. Comput. Appl. Math. 260 (2014), pp. 244–257.
- [5] A. Beck and Y. Vaisbourd, *Globally solving the trust region subproblem using simple first-order methods*, SIAM J. Optim. 28(3) (2018), pp. 1951–1967.
- [6] T. Bianconcini, G. Liuzzi, B. Morini, and M. Sciandrone, *On the use of iterative methods in cubic regularization for unconstrained optimization*, Comput. Optim. Appl. 60(1) (2015), pp. 35–57.
- [7] Y. Carmon and J.C. Duchi, *Gradient descent efficiently finds the cubic-regularized non-convex Newton step*, (2016). Available at arXiv:1612.00547.
- [8] Y. Carmon and J.C. Duchi, *Analysis of Krylov subspace solutions of regularized nonconvex quadratic problems*, (2018). Available at arXiv:1806.09222v1.
- [9] C. Cartis, N.I.M. Gould, and Ph. L. Toint, *Adaptive cubic regularisation methods for unconstrained optimization. Part I: Motivation, convergence and numerical results*, Math. Program. Ser. A 127(2) (2011), pp. 245–295.

- [10] C. Cartis, N.I.M. Gould, and M. Lange, *On monotonic estimates of the norm of the minimizers of regularized quadratic functions in Krylov spaces*. Tech. Rep. RAL-TR-2019, Rutherford Appleton Laboratory, Chilton, Oxfordshire, England, 2019.
- [11] A.R. Conn, N.I.M. Gould, and Ph. L. Toint, *Trust-Region Methods*, SIAM, Philadelphia, 2000.
- [12] S. Demko, W.F. Moss, and P.W. Smith, *Decay rates for inverses of band matrices*, *Math. Comput.* 43(168) (1984), pp. 491–499.
- [13] J.B. Erway and P.E. Gill, *A subspace minimization method for the trust-region step*, *SIAM J. Optim.* 20(3) (2009), pp. 1439–1461.
- [14] J.B. Erway, P.E. Gill, and J.D. Griffin, *Iterative methods for finding a trust-region step*, *SIAM J. Optim.* 20(2) (2009), pp. 1110–1131.
- [15] C. Fortin and H. Wolkowicz, *The trust region subproblem and semidefinite programming*, *Optim. Methods Softw.* 19(1) (2004), pp. 41–67.
- [16] D.M. Gay, *Computing optimal locally constrained steps*, *SIAM J. Sci. Statist. Comput.* 2(2) (1981), pp. 186–197.
- [17] G.H. Golub and C.F. Van Loan, *Matrix Computations*, 3rd ed., Johns Hopkins University Press, Baltimore, 1996.
- [18] N.I.M. Gould, S. Lucidi, M. Roma, and Ph. L. Toint, *Solving the trust-region subproblem using the Lanczos method*, *SIAM J. Optim.* 9(2) (1999), pp. 504–525.
- [19] N.I.M. Gould, D. Orban, and Ph. L. Toint, *GALAHAD—a library of thread-safe fortran 90 packages for large-scale nonlinear optimization*, *ACM Trans. Math. Soft.* 29(4) (2003), pp. 353–372.
- [20] N.I.M. Gould, D.P. Robinson, and H.S. Thorne, *On solving trust-region and other regularised subproblems in optimization*, *Math. Program. Comput.* 2(1) (2010), pp. 21–57.
- [21] N.I.M. Gould, D. Orban, and Ph. L. Toint, *CUTEst: A constrained and unconstrained testing environment with safe threads for mathematical optimization*, *Comput. Optim. Appl.* 60(3) (2015), pp. 545–557.
- [22] A. Greenbaum, *Iterative Methods for Solving Linear Systems*, SIAM, Philadelphia, 1997.
- [23] W.W. Hager, *Minimizing a quadratic over a sphere*, *SIAM J. Optim.* 12(1) (2001), pp. 188–208.
- [24] E. Hazan and T. Koren, *A linear-time algorithm for trust region problems*, *Math. Program.* 158(1–2) (2016), pp. 363–381.
- [25] N. Ho-Nguyen and F. Kilinç-Karzan, *A second-order cone based approach for solving the trust-region subproblem and its variants*, (2016). Available at arXiv:1603.03366.
- [26] J. Lampe, M. Rojas, D.C. Sorensen, and H. Voss, *Accelerating the LSTRS algorithm*, *SIAM J. Sci. Comput.* 33(1) (2011), pp. 175–194.
- [27] J. Liesen and Z. Strakoš, *Krylov Subspace Methods*, Oxford University Press, Oxford, 2013.
- [28] L. Lukšan, C. Matonoha, and J. Vlček, *On Lagrange multipliers of trust-region subproblems*, *BIT Numer. Math.* 48(4) (2008), pp. 763–768.
- [29] J.M. Martínez, *Local minimizers of quadratic functions on Euclidean balls and spheres*, *SIAM J. Optim.* 4(1) (1994), pp. 159–176.
- [30] J.J. Moré and D.C. Sorensen, *Computing a trust region step*, *SIAM J. Sci. Statist. Comput.* 4(3) (1983), pp. 553–572.
- [31] Yu. Nesterov and B.T. Polyak, *Cubic regularization of Newton method and its global performance*, *Math. Program.* 108(1) (2006), pp. 177–205.
- [32] J. Nocedal and S.J. Wright, *Numerical Optimization*, 2nd ed., Series in Operations Research. Springer Verlag, Heidelberg, Berlin, New York, 2006.
- [33] F. Rendl and H. Wolkowicz, *A semidefinite framework for trust region subproblems with applications to large scale minimization*, *Math. Program. Ser. B* 77(2) (1997), pp. 273–299.
- [34] M. Rojas, S.A. Santos, and D.C. Sorensen, *A new matrix-free algorithm for the large-scale trust-region subproblem*, *SIAM J. Optim.* 11(3) (2001), pp. 611–646.
- [35] M. Rojas, S.A. Santos, and D.C. Sorensen, *Algorithm 873: LSTRS: MATLAB software for large-scale trust-region subproblems and regularization*, *ACM Trans. Math. Softw.* 34(2) (2008), Article 11.
- [36] C.W. Royer and S.J. Wright, *Complexity analysis of second-order line-search algorithms for smooth nonconvex optimization*, (2017). Available at arXiv:1706.03131v2.

- [37] T. Steihaug, *The conjugate gradient method and trust regions in large scale optimization*, SIAM J. Numer. Anal. 20(3) (1983), pp. 626–637.
- [38] Ph. L. Toint, *Towards an efficient sparsity exploiting Newton method for minimization*, in *Sparse Matrices and Their Uses*, I.S. Duff, eds., Academic Press, London, 1981, pp. 57–88.
- [39] H.A. van der Vorst, *Iterative Krylov Methods for Large Linear Systems*, Cambridge University Press, Cambridge, 2003.
- [40] J. Wong and Y. Xia, *A linear-time algorithm for the trust region subproblem based on hidden convexity*, Optim. Lett. 11(8) (2017), pp. 1639–1646.
- [41] L.-H. Zhang and C. Shen, *A nested Lanczos method for the trust-region subproblem*, SIAM J. Sci. Comput. 40(4) (2018), pp. A2005–A2032.
- [42] L.-H. Zhang, C. Shen, and R.-C. Li, *On the generalized Lanczos trust-region method*, SIAM J. Optim. 27(3) (2017), pp. 2110–2142.

Appendices

Appendix 1

Using (50) to compare the model decrease from x_k to x_m , we deduce that

$$\begin{aligned}
 q(x_k) - q(x_m) &= \frac{1}{2}x_k^T Hx_k + g^T x_k - \frac{1}{2}x_m^T Hx_m - g^T x_m \\
 &= \frac{1}{2}x_k^T Hx_k - \frac{1}{2}x_m^T Hx_m + g^T (x_k - x_m) \\
 &= \frac{1}{2}x_k^T Hx_k - \frac{1}{2}x_m^T Hx_m + (x_m - x_k)^T Hx_m + \mu_m (x_m - x_k)^T x_m \\
 &= \frac{1}{2}(x_m - x_k)^T (H + \mu_m I)(x_m - x_k) + \frac{1}{2}\mu_m (\|x_m\|^2 - \|x_k\|^2). \tag{A1}
 \end{aligned}$$

Since (15) shows that $H + \mu_m I$ is positive definite on \mathcal{K}_m , and as $r_k \in \mathcal{K}_m$, it follows that

$$(H + \mu_m I)^{-1} r_k = (x_k - x_m) + (\mu_k - \mu_m)(H + \mu_m I)^{-1} x_k,$$

and hence, taking the inner product with r_k from (51),

$$\begin{aligned}
 r_k^T (H + \mu_m I)^{-1} r_k &= (x_m - x_k)^T (H + \mu_m I)(x_m - x_k) \\
 &\quad + 2(\mu_m - \mu_k)(x_m - x_k)^T x_k + (\mu_m - \mu_k)^2 x_k^T (H + \mu_m I)^{-1} x_k \\
 &= 2[q(x_k) - q(x_m)] - \mu_m [\|x_m\|^2 - \|x_k\|^2] \\
 &\quad + 2(\mu_m - \mu_k)(x_m - x_k)^T x_k + (\mu_m - \mu_k)^2 x_k^T (H + \mu_m I)^{-1} x_k \\
 &= 2[q(x_k) - q(x_m)] + \mu_k (\|x_k\|^2 - \|x_m\|^2) \\
 &\quad - (\mu_m - \mu_k) \|x_m - x_k\|^2 + (\mu_m - \mu_k)^2 x_k^T (H + \mu_m I)^{-1} x_k \tag{A2}
 \end{aligned}$$

using (A1). As $r_k \in \mathcal{K}_m$, referring back to (7), we have that

$$r_k = U_+ \hat{r}_k = U \begin{pmatrix} \hat{r}_k \\ 0 \end{pmatrix},$$

for some \hat{r}_k , and thus

$$r_k^T (H + \mu_m I)^{-1} r_k = \hat{r}_k^T (\Lambda_+ + \mu_m I)^{-1} \hat{r}_k \geq \frac{\|\hat{r}_k\|^2}{(\lambda_+^{\max} + \mu_m)} = \frac{\|r_k\|^2}{(\lambda_+^{\max} + \mu_m)}$$

where we recall the definition of λ_+^{\max} from (21). Hence (A2) leads directly to (52).

We recall that [10,37]

$$\|x_k\| \leq \|x_m\| \tag{A3}$$

and [10,28]

$$0 \leq \mu_k \leq \mu_m, \tag{A4}$$

and hence $\mu_k (\|x_k\|^2 - \|x_m\|^2) \leq 0$.

Suppose that $\mu_m = 0$. Then (A4) implies that $\mu_k = 0$, while (50) gives $Hx_m = -g$ and thus (1), (7), (8) and (21) combine to give

$$q(0) - q(x_m) = \frac{1}{2} \bar{g}_+^T \Lambda_+^{-1} \bar{g}_+ \leq \frac{\|\bar{g}_+\|^2}{(\lambda_+^{\min} + \mu_m)} = \frac{\|g\|^2}{(\lambda_+^{\min} + \mu_m)}. \quad (\text{A5})$$

Combining (22) (52) and (A5) gives

$$\|r_k\|^2 \leq 2(\lambda_+^{\max} + \mu_m)[q(x_k) - q(x_m)] \leq 72\kappa_m \left(e^{-4/\sqrt{\kappa_m}} \right)^k \|g\|^2,$$

i.e.

$$\|r_k\| \leq 6\sqrt{2}\sqrt{\kappa_m} \left(e^{-2/\sqrt{\kappa_m}} \right)^k \|g\|. \quad (\text{A6})$$

We note that obtaining an error bound via (50)–(A2) is similar to the approach taken by [36, §3.2] in the absence of a trust region.

Unfortunately, it is unclear how to proceed when $\mu_m > 0$ —there are two sub-cases $\mu_k = 0$ and $\mu_k > 0$, but we cannot see a way for either. The issue, of course, is the extra term

$$-(\mu_m - \mu_k)\|x_m - x_k\|^2 + (\mu_m - \mu_k)^2 x_k^T (H + \mu_m I)^{-1} x_k \quad (\text{A7})$$

in (52). Ideally, we would like to show that this is negative in which case a bound of the form (A6) would follow. We also have a bound

$$\begin{aligned} x_k^T (H + \mu_m I)^{-1} x_k &= (x_k - x_m)^T (H + \mu_m I)^{-1} (x_k - x_m) \\ &\quad + 2x_k^T (H + \mu_m I)^{-1} x_m - x_m^T (H + \mu_m I)^{-1} x_m \\ &\leq (x_k - x_m)^T (H + \mu_m I)^{-1} (x_k - x_m) - 2g^T x_k \end{aligned}$$

on the second term in (A7), but that doesn't seem to help. Another possibility is to show that the two terms in (A7) decay exponentially, although we see no reason why in particular $(\mu_m - \mu_k)$ would—one might for example have $\mu_k = 0$ for all $k = 1, \dots, m-1$ (i.e. the trust-region constraint is inactive), but $\mu_m > 0$ (the constraint becomes active).

Appendix 2

Our second attempt to find a useful bound on the residual is based on the relationships (26)–(27), and uses the following identity.

Lemma A.1: *For any scalar λ , we have*

$$V_k^T (H + \lambda I)^j V_k e_1 = (T_k + \lambda I)^j e_1 \quad (\text{A8})$$

for $j = 1, \dots, k$.

Proof: We first show that

$$(H + \lambda I)^j V_k = V_k (T_k + \lambda I)^j + \gamma_k \sum_{i=0}^{j-1} (H + \lambda I)^{j-i-1} v_{k+1} e_k^T (T_k + \lambda I)^i \quad (\text{A9})$$

for all $k \geq 1$. This follows immediately when $j = 1$ as

$$(H + \lambda I)V_k = V_k(T_k + \lambda I) + \gamma_k v_{k+1} e_k^T \quad (\text{A10})$$

from (26). Suppose that (A9) holds for some $j > 1$. Then multiplying by $H + \lambda I$ and using (A10), we have

$$\begin{aligned}
 (H + \lambda I)^{j+1} V_k &= (H + \lambda I) V_k (T_k + \lambda I)^j + \gamma_k \sum_{i=0}^{j-1} (H + \lambda I)^{j-i} v_{k+1} e_k^T (T_k + \lambda I)^i \\
 &= (V_k (T_k + \lambda I) + \gamma_k v_{k+1} e_k^T) (T_k + \lambda I)^j + \gamma_k \sum_{i=0}^{j-1} (H + \lambda I)^{j-i} v_{k+1} e_k^T (T_k + \lambda I)^i \\
 &= V_k (T_k + \lambda I)^{j+1} + \gamma_k v_{k+1} e_k^T (T_k + \lambda I)^j + \gamma_k \sum_{i=0}^{j-1} (H + \lambda I)^{j-i} v_{k+1} e_k^T (T_k + \lambda I)^i \\
 &= V_k (T_k + \lambda I)^{j+1} + \gamma_k \sum_{i=0}^j (H + \lambda I)^{j-i} v_{k+1} e_k^T (T_k + \lambda I)^i,
 \end{aligned}$$

and thus (A9) holds for $j + 1$. Hence (A9) holds for all $j \geq 1$ by induction.

Now observe that $(T_k + \lambda I)e_1$ only has nonzeros in positions 1 and 2, $(T_k + \lambda I)^2 e_1 = (T_k + \lambda I)((T_k + \lambda I)e_1)$ only has nonzeros in positions 1 to 3, and in general $(T_k + \lambda I)^{j-1} e_1 = (T_k + \lambda I)((T_k + \lambda I)^{j-2} e_1)$ only has nonzeros in positions 1 to j . Thus from (A9) and the orthogonality of the columns of V_k , we have

$$\begin{aligned}
 V_k^T (H + \lambda I)^j V_k e_1 &= V_k^T V_k (T_k + \lambda I)^j e_1 + \gamma_k \sum_{i=0}^{j-1} V_k^T (H + \lambda I)^{j-i-1} v_{k+1} e_k^T (T_k + \lambda I)^i e_1 \\
 &= (T_k + \lambda I)^j e_1 + \gamma_k \sum_{i=0}^{j-2} V_k^T (H + \lambda I)^{j-i-1} v_{k+1} e_k^T (T_k + \lambda I)^i e_1 + \gamma_k V_k^T v_{k+1} e_k^T (T_k + \lambda I)^{j-1} e_1 \\
 &= (T_k + \lambda I)^j e_1
 \end{aligned}$$

as required, since $e_k^T (T_k + \lambda I)^i e_1 = 0$ for $i = 0, \dots, j-2$ and $j = 1, \dots, k$, and $V_k^T v_{k+1} = 0$. ■

Since $x_k \in \mathcal{K}_k = \text{span}\{H^i g\}_{i=0}^{k-1} \equiv \text{span}\{(H + \mu_k I)^i g\}_{i=0}^{k-1}$, we have

$$x_k = \sum_{j=0}^{k-1} \eta_j (H + \mu_k I)^j g$$

for coefficients $\eta_j, j = 0, \dots, k-1$, and thus

$$\begin{aligned}
 r_k &= g + (H + \mu_k I) \sum_{j=0}^{k-1} \eta_j (H + \mu_k I)^j g = g + \sum_{j=0}^{k-1} \eta_j (H + \mu_k I)^{j+1} g \\
 &= g + \sum_{j=1}^k \eta_{j-1} (H + \mu_k I)^j g \\
 &= \psi_k (H + \mu_k I) g \tag{A11}
 \end{aligned}$$

for some

$$\psi_k(\lambda) = \sum_{j=0}^k \omega_j \lambda^j \in \mathcal{P}_k = \{\text{polynomials } \psi \text{ of degree } k \text{ for which } \psi(0) = 1\}.$$

But then

$$\begin{aligned}
 V_k^T r_k &= V_k^T \psi_k(H + \mu_k I)g = \sum_{j=0}^k \omega_j V_k^T (H + \mu_k I)^j g \\
 &= \|g\| \sum_{j=0}^k \omega_j V_k^T (H + \mu_k I)^j V_k e_1 = \|g\| \sum_{j=0}^k \omega_j (T_k + \mu_k I)^j e_1 \\
 &= \|g\| \psi_k(T_k + \mu_k I) e_1
 \end{aligned} \tag{A12}$$

using (30) and (A8). Since T_k is irreducible, it has distinct eigenvalues $\theta_{i,k}$, $i = 1, \dots, k$, and as (31) indicates that $V_k^T r_k = 0$, (A12) implies that ψ_k is a scalar multiple of the minimum polynomial

$$\phi_k(\lambda) = \prod_{i=1}^k (\lambda - \theta_{i,k} - \mu_k)$$

of the irreducible matrix $T_k + \mu_k I$. Indeed, $\psi_k(\lambda) = \phi_k(\lambda)/\phi_k(0)$ since we require that $\psi_k(0) = 1$.

Referring back to (7) and (8), we have that $H = U \Lambda U^T$ and $g = U \bar{g}$ for matrices U of eigenvectors and Λ of eigenvalues. Then (A11) gives

$$r_k = \psi_k \left(U(\Lambda + \mu_k I)U^T \right) U \bar{g} = U \psi_k(\Lambda + \mu_k I) \bar{g} = U_+ \psi_k(\Lambda_+ + \mu_k I) \bar{g}_+,$$

and hence

$$\begin{aligned}
 \|r_k\|^2 &= \|\psi_k(\Lambda_+ + \mu_k I) \bar{g}_+\|^2 = \sum_{j \in \mathcal{I}_+} \psi_k^2(\lambda_j + \mu_k) \bar{g}_j^2 \\
 &\leq \max_{j \in \mathcal{I}_+} \psi_k^2(\lambda_j + \mu_k) \sum_{j \in \mathcal{I}_+} \bar{g}_j^2 = \max_{j \in \mathcal{I}_+} \psi_k^2(\lambda_j + \mu_k) \|\bar{g}_*\|^2 \\
 &= \max_{j \in \mathcal{I}_+} \psi_k^2(\lambda_j + \mu_k) \|g\|^2.
 \end{aligned}$$

This then provides the estimate (53).