

Solving the cubic regularization model by a nested restarting Lanczos method

Xiaojing Jia* Xin Liang† Chungen Shen‡ Lei-Hong Zhang§

Abstract

As a variant of the classical trust-region method for unconstrained optimization, the cubic regularization of Newton method introduces a cubic regularization term in the surrogate objective to adaptively adjust the updating step and deals with cases with both indefinite and definite Hessians. It has been demonstrated that the cubic regularization of Newton method enjoys a good global convergence and is an efficient solver for the unconstrained minimization. The main computational cost in each iteration is to solve a cubic regularization subproblem. The Newton iteration is a common and efficient method for this task, especially for small-to medium-size problems. For large size problems, a Lanczos type method was proposed in [Cartis, Gould and Toint, Math. Program., 127:245–295(2011)]. This method relies on a Lanczos procedure to reduce the large-scale cubic regularization subproblem to a small one and solve it by the Newton iteration. For large and ill-conditioned problems, the Lanczos method still needs to produce a large dimensional subspace to achieve a relatively highly accurate approximation, which declines its performance overall. In this paper, we first show that the cubic regularization subproblem can be equivalently transformed into a quadratic eigenvalue problem, which provides an eigensolver alternative to the Newton iteration. We then establish the convergence of the Lanczos method and also propose a nested restarting version for the large scale and ill-conditioned case. By integrating the nested restarting Lanczos iteration into the cubic regularization of Newton method, we verify its efficiency for solving large scale minimization problems in CUTEst collection.

Key words. Unconstrained optimization, Cubic regularization, Newton’s method, Eigenvalue problem, Global convergence, Lanczos process, Restarting

AMS subject classifications. 90C30, 90C06, 65K05, 49M15, 65F15

1 Introduction

For the unconstrained minimization

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}), \quad (1.1)$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and its Hessian $\nabla^2 f(\mathbf{x}) \in \mathbb{R}^{n \times n}$ satisfies

$$\|\nabla^2 f(\mathbf{x}) - \nabla^2 f(\mathbf{y})\|_2 \leq L\|\mathbf{x} - \mathbf{y}\|_2, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n,$$

*School of Mathematics, Shanghai University of Finance and Economics, 777 Guoding Road, Shanghai 200433, China. Email: 13127761671@163.com.

†Yau Mathematical Sciences Center, Tsinghua University, Beijing 100084, China, and Yanqi Lake Beijing Institute of Mathematical Sciences and Applications, Beijing 101408, China. Supported in part by the National Natural Science Foundation of China NSFC-11901340 and NSFC-12071332. Email: liangxinslm@tsinghua.edu.cn.

‡College of Science, University of Shanghai for Science and Technology, Shanghai 200093, China. Email: shenchungen@gmail.com.

§Corresponding author. School of Mathematical Sciences, Soochow University, Suzhou 215006, Jiangsu, China. Supported in part by the National Natural Science Foundation of China NSFC-12071332. Email: longzlh@suda.edu.cn.

[27] proposes a variant of the classical trust-region method [9, 35], namely, a cubic regularization of Newton method. Let $0 < L_0 \leq L$. At the current iterate \mathbf{x}_k , the method solves the following cubic regularization model:

$$\min_{\mathbf{h} \in \mathbb{R}^n} \left\{ m(\mathbf{h}) := f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^\top \mathbf{h} + \frac{1}{2} \mathbf{h}^\top \nabla^2 f(\mathbf{x}_k) \mathbf{h} + \frac{M_k}{6} \|\mathbf{h}\|_2^3 \right\}, \quad (1.2)$$

with a properly chosen $M_k \in [L_0, 2L]$; whenever $f(\mathbf{x}_k + \mathbf{h}) \leq m(\mathbf{h})$, the iterate \mathbf{x}_k is updated to $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{h}$. The global convergence and local quadratic convergence of this iteration are proved in [27]. Moreover, efficient variants and modifications have been extensively discussed in [7, 8]. In particular, [7, 8] show that there is flexibility in using certain symmetric approximates $H_k \in \mathbb{R}^{n \times n}$ of the Hessian matrix $\nabla^2 f(\mathbf{x}_k)$ where the nice global and local convergence can still be guaranteed. Moreover, it is shown in [8] that an adaptive cubic regularization of Newton method needs at most $O(\epsilon^{-3/2})$ function- and gradient-evaluations to achieve an approximation \mathbf{x}_k with $\|\nabla f(\mathbf{x}_k)\|_2 \leq \epsilon$ for a given accuracy ϵ . This improves the worst-case complexity [19] of the traditional second-order trust-region method where $O(\epsilon^{-2})$ iterations are required to have $\|\nabla f(\mathbf{x}_k)\|_2 \leq \epsilon$.

The above cubic regularization of Newton method, as was claimed in [7, 8], can be viewed as an adaptive version of the classical trust-region method [9], where the rules for updating M_k are justified by an analogy to trust-region methods; in particular, M_k might be regarded as the reciprocal of the trust-region radius. Instead of imposing a trust-region $\|\mathbf{h}\| \leq \Delta$ for the well-definiteness when $\nabla^2 f(\mathbf{x}_k)$ is indefinite, the cubic regularization model (1.2) introduces a regularization term $\frac{M_k}{3} \|\mathbf{h}\|_2 I_n$ to $\nabla^2 f(\mathbf{x}_k)$ to adaptively adjust the solution \mathbf{h} both for the indefinite case and for the definite case. It can be seen that when $\nabla^2 f(\mathbf{x}_k)$ is indefinite, the solution \mathbf{h} cannot be of infinite norm as the regularization term $\frac{M_k}{3} \|\mathbf{h}\|_2 I_n$ will enforce the modified Hessian $\nabla^2 f(\mathbf{x}_k) + \frac{M_k}{3} \|\mathbf{h}\|_2 I_n$ to be positive definite, and hence the solution of (2.1) is well-defined.

Regarding the update $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{h}$ in the cubic regularization of Newton method, we note that the main computational step lies in solving the cubic regularization model (1.2). This is similar to the trust-region method where the primary computation is to solve the so-called trust-region subproblem (TRS)

$$\min_{\|\mathbf{h}\|_W \leq \Delta} \mathbf{h}^\top \nabla f(\mathbf{x}_k) + \frac{1}{2} \mathbf{h}^\top \nabla^2 f(\mathbf{x}_k) \mathbf{h} \quad (1.3)$$

where W is a proper positive definite weighted matrix. The TRS has been well-understood in theory (see e.g., [9, 21, 22, 26, 28, 38]) and many efficient numerical methods have been proposed which can be basically grouped into factorization-based algorithms for small-to-medium sized dense problems (see, e.g., [1, 22, 23, 26]) and factorization-free algorithms for large scale and sparse problems (see, e.g., [14, 17, 21, 30, 31, 32, 33, 34, 35, 37, 40, 41]). One of the most widely known factorization-based methods is the Moré-Sorensen method [26], which is a Newton method for solving the associated Lagrange multiplier. Another approach proposed recently in [1] generalizes [14] and translates TRS into certain eigenvalue problems. For large scale TRS or when the Hessian matrix $\nabla^2 f(\mathbf{x}_k)$ is only available through its action $\nabla^2 f(\mathbf{x}_k) \mathbf{z}$ on a vector \mathbf{z} , a Krylov subspace method, namely the generalized Lanczos Trust-Region method (GLTR), was proposed by Gould *et al.* [15] (see also [9, Chapter 5]). The convergence of GLTR has been recently established in [5, 6, 18, 40, 41] and reveals the linear convergence in the worst case scenario [41, 5].

The purpose of this paper is to develop efficient methods for (1.2). A recent work by Lieder [25] extends [1] for TRS and proposes an equivalent $2(n+1)$ dimensional generalized eigenvalue problem to solve (1.2). In this paper, we shall first introduce an equivalent $(n+1)$ dimensional quadratic eigenvalue problem (QEP) for (1.2). By the new QEP, on the one hand, efficient Krylov subspace methods working on \mathbb{R}^{n+1} such as the second-order Arnoldi process (SOAR) [2] can be directly applied to solve (1.2), and provides new forms of $2(n+1)$ dimensional generalized eigenvalue problems for (1.2) on the other hand. The equivalent reformulations in the form of generalized eigenvalue problem and QEP provide relations of (1.2) with the eigenvalue problem and also offer numerical schemes to solve (1.2), especially for small to- medium size cases. For large size problems or the cases when only the action $\nabla^2 f(\mathbf{x}_k) \mathbf{z}$ of $\nabla^2 f(\mathbf{x}_k)$ on a vector \mathbf{z} is available, we discuss Lanczos methods [7, Section 6.2] for (1.2). Previously, the linear convergence of this Lanczos method has been established in [5, 6], and we will sharpen this convergence result

81 by developing a new convergence analysis; furthermore, we will also design an efficient restarting
82 scheme for this Lanczos method to obtain accurate approximation for **ill-conditioned instances**¹
83 of (1.2). The resulting approach consists of a nested restarting procedure and is able to alleviate
84 numerical difficulties of the basic Lanczos method [5, 6] caused by the dimension increment of the
85 underlying Krylov subspace in the Lanczos process. As a practical application, we will integrate
86 the nested restarting Lanczos method for (1.2) to solve large scale minimization problems (1.1)
87 from CUTEst collection [16]. Our numerical experience demonstrates that the nested restarting
88 Lanczos method can be an efficient approach to deal with ill-conditioned inner subproblems (1.2)
89 and improves the overall performance of the cubic regularization of Newton method.

90 We organize the paper in the following way: in section 2, we first provide basic properties of
91 (1.2). Section 3 then introduces a QEP and establishes the equivalence. The presentation of the
92 Lanczos method of (1.2) is given in section 4 where we shall establish the linear convergence in the
93 worst case, and also propose a nested restarting version for ill-conditioned problems. Numerical
94 verification of our restarting Lanczos method will be carried out in section 5, and final conclusions
95 are drawn in section 6.

96 **Notation.** We use the following notation system in this paper. Vectors are generally referred
97 to be column vectors and are typeset in bold lower case letters; in particular, $\mathbf{e}_i \in \mathbb{R}^n$ is the i th
98 column of the identity matrix I_n . For a matrix $A \in \mathbb{R}^{m \times n}$, its Moore-Penrose inverse and its
99 column range space of A are presented by A^\dagger and $\text{span}(A)$, respectively. The dimension of $\text{span}(A)$
100 is given by $\dim(\text{span}(A))$. To facilitate the presentation, we shall conveniently adopt the MATLAB
101 format to access the entries of vectors and matrices: $A_{(i,j)}$ is (i, j) th entry of A , and $A_{(k:\ell, i:j)}$ is
102 the submatrix of A that contains intersections of row k to row ℓ and column i to column j . For a
103 square matrix A , the set of all eigenvalues and the determinant are denoted by $\text{eig}(A)$ and $\det(A)$,
104 respectively. Finally, $\mathcal{K}_\ell(A, \mathbf{x})$ stands for the ℓ th Krylov subspace and any $\mathbf{h} \in \mathcal{K}_\ell(A, \mathbf{x})$ can be
105 expressed by $\mathbf{h} = p(A)\mathbf{x}$, where $p \in \mathbb{P}_\ell$ is a polynomial with degree no higher than ℓ .

106 2 The cubic regularization model

For the simplicity of presentation, we omit the subscript k in (1.2), and also denote

$$\mathbf{g} = \nabla f(\mathbf{x}_k), \quad \sigma = M_k/2, \quad H = H_k \approx \nabla^2 f(\mathbf{x}_k).$$

107 Thus, we focus on the following minimization:

$$108 \quad \min_{\mathbf{h} \in \mathbb{R}^n} \left\{ m(\mathbf{h}) := \mathbf{g}^\top \mathbf{h} + \frac{1}{2} \mathbf{h}^\top H \mathbf{h} + \frac{\sigma}{3} \|\mathbf{h}\|_2^3 \right\}. \quad (2.1)$$

109 The following result generalizes the well-known sufficient and necessary conditions (Gay [13] and
110 Moré and Sorensen [26]) for the trust-region subproblem to (2.1).

111 **Theorem 2.1.** ([7, Theorem 3.1] and [27, Theorem 10]) *Any \mathbf{h}_{opt} is a global minimizer of (2.1)*
112 *over \mathbb{R}^n if and only if it satisfies the system of equations*

$$113 \quad (H + \lambda_{\text{opt}} I_n) \mathbf{h}_{\text{opt}} = -\mathbf{g}, \quad (2.2)$$

114 *where $\lambda_{\text{opt}} = \sigma \|\mathbf{h}_{\text{opt}}\|_2$ and $H + \lambda_{\text{opt}} I_n$ is positive semidefinite. If $H + \lambda_{\text{opt}} I_n$ is positive definite,*
115 *then \mathbf{h}_{opt} is unique.*

116 By Theorem 2.1, we know that $\lambda_{\text{opt}} \geq -\theta_1^+$ where

$$117 \quad -\theta_1^+ := \max(0, -\theta_1) \quad \text{and} \quad \theta_1 = \theta_2 = \cdots = \theta_p < \theta_{p+1} \leq \cdots \leq \theta_n \quad (2.3)$$

118 are the ordered eigenvalues² of H and $H = U\Theta U^\top$ is its spectral decomposition with $U =$
119 $[\mathbf{u}_1, \dots, \mathbf{u}_n]$ orthonormal. Notice that the value λ_{opt} plays a similar role with the Lagrangian

¹A problem (1.2) is said to be ill-conditioned if the matrix $\nabla^2 f(\mathbf{x}_k) + \frac{M_k}{3} \|\mathbf{h}_{\text{opt}}\|_2 I_n$ is ill-conditioned, where \mathbf{h}_{opt} is the minimizer of (1.2); see Theorem 4.1.

²In our discussion, without loss of generality, we assume that $p < n$. When $p = n$, then Theorem 2.1 implies that \mathbf{h}_{opt} is parallel to \mathbf{g} , and \mathbf{h}_{opt} can be obtained by solving a one-dimensional minimization.

120 multiplier for the trust-region subproblem; thus, in what follows, we will also call λ_{opt} as the La-
 121 grangian multiplier for the problem (2.1). Now, by the fact $\lambda_{\text{opt}} = \sigma \|\mathbf{h}_{\text{opt}}\|_2$, whenever $\lambda_{\text{opt}} > -\theta_1$,
 122 λ_{opt} can be found via the system:

$$123 \quad \lambda_{\text{opt}} = \sigma \|\mathbf{h}_{\text{opt}}\|_2 = \sigma \|(H + \lambda_{\text{opt}} I_n)^{-1} \mathbf{g}\|_2,$$

124 OR

$$125 \quad \lambda_{\text{opt}}^2 = \sigma^2 \mathbf{g}^T (H + \lambda_{\text{opt}} I_n)^{-2} \mathbf{g} = \mathbf{t}^T (\Theta + \lambda_{\text{opt}} I_n)^{-2} \mathbf{t} = \sum_{j=1}^n \frac{t_j^2}{(\theta_j + \lambda_{\text{opt}})^2}, \quad (2.4)$$

126 where $\mathbf{t} = \sigma U^T \mathbf{g} = [t_1, \dots, t_n]^T$. Introduce the system

$$127 \quad q(\lambda) := \lambda^2 - \sum_{j=1}^n \frac{t_j^2}{(\theta_j + \lambda)^2}, \quad (2.5)$$

128 and we have $q(\lambda_{\text{opt}}) = 0$.

129 We next present results for special cases: $\mathbf{g} = \mathbf{0}$ and $\lambda_{\text{opt}} = \max(0, -\theta_1)$.

130 **Theorem 2.2.** For (2.1) with $\sigma > 0$, we have

131 (i) if $\mathbf{g} = \mathbf{0}$, then $\lambda_{\text{opt}} = -\theta_1^+$;

132 (ii) $\lambda_{\text{opt}} = 0$ if and only if $\mathbf{g} = \mathbf{0}$ and H is positive semidefinite;

133 (iii) $\lambda_{\text{opt}} = -\theta_1$ if and only if

$$134 \quad \mathbf{g} \perp \mathcal{E}_1 = \text{span}(\{\mathbf{u}_1, \dots, \mathbf{u}_p\}) \quad \text{and} \quad -\theta_1 \geq \sigma \|(H - \theta_1 I_n)^\dagger \mathbf{g}\|_2, \quad (2.6)$$

135 where \mathcal{E}_1 is the eigenspace associated with the smallest eigenvalue $\theta_1 = \dots = \theta_p$ of H .

136 *Proof.* For (i), we know that when $-\theta_1 \leq 0$, then by (2.2), the assumption $\lambda_{\text{opt}} > 0$ leads to
 137 $\mathbf{h}_{\text{opt}} = \mathbf{0}$, which contradicts with $0 < \lambda_{\text{opt}} = \sigma \|\mathbf{h}_{\text{opt}}\|_2 = 0$. Therefore, when $-\theta_1 \leq 0$, we have
 138 $\lambda_{\text{opt}} = \max(0, -\theta_1) = 0$. For $-\theta_1 > 0$, by a similar argument, we conclude $\lambda_{\text{opt}} = -\theta_1^+ = -\theta_1$.

139 For (ii), if $\lambda_{\text{opt}} = 0$, by $0 = \lambda_{\text{opt}} = \sigma \|\mathbf{h}_{\text{opt}}\|_2$ and the semidefiniteness of $H + \lambda_{\text{opt}} I_n$, we have
 140 $\mathbf{h}_{\text{opt}} = \mathbf{0}$, and $\theta_1 \geq 0$; therefore, $\mathbf{0} = H \mathbf{h}_{\text{opt}} = -\mathbf{g}$. The converse is from (i).

For (iii), we first consider the necessity for $\lambda_{\text{opt}} = -\theta_1$. If $\theta_1 = 0$, then by (ii), we know that
 (2.6) holds. If $\lambda_{\text{opt}} = -\theta_1 > 0$, then optimality condition (2.2) implies that $\mathbf{g} \perp \mathcal{E}_1$, and also all the
 solutions to the system (2.2) can be given by $\mathbf{h} = -(H - \theta_1 I_n)^\dagger \mathbf{g} + \mathbf{u}$ where $\mathbf{u} \in \mathcal{E}_1$ is arbitrary.
 Also, the condition $-\theta_1 = \lambda_{\text{opt}} = \sigma \|\mathbf{h}_{\text{opt}}\|_2$ gives

$$\theta_1^2 = \sigma^2 [\|(H - \theta_1 I_n)^\dagger \mathbf{g}\|_2^2 + \|\mathbf{u}\|_2^2] = \sum_{i=p+1}^n \frac{t_i^2}{(\theta_i - \theta_1)^2} + \sigma^2 \|\mathbf{u}\|_2^2, \quad \text{with } \mathbf{t} = \sigma U^T \mathbf{g} = [t_1, \dots, t_n]^T.$$

141 Consider the function

$$142 \quad \tilde{q}(\lambda) = \lambda^2 - \sum_{i=p+1}^n \frac{t_i^2}{(\theta_i + \lambda)^2}. \quad (2.7)$$

143 It is easy to see that $\tilde{q}(\lambda) \rightarrow +\infty$ as $\lambda \rightarrow +\infty$, and $\tilde{q}(\lambda) > 0$ for $\lambda \in [-\theta_1, +\infty)$. Thus,
 144 $\tilde{q}(-\theta_1) - \sigma^2 \|\mathbf{u}\|_2^2 = 0$ for some \mathbf{u} holds only if $-\theta_1 \geq \sigma \|(H - \theta_1 I_n)^\dagger \mathbf{g}\|_2$, and in this case, furthermore,
 145 the optimal solution $\mathbf{h}_{\text{opt}} = -(H - \theta_1 I_n)^\dagger \mathbf{g} + \mathbf{u}$ where \mathbf{u} is any eigenvector of H corresponding to
 146 θ_1 with $\|\mathbf{u}\|_2^2 = \theta_1^2 / \sigma^2 - \|(H + \lambda I_n)^\dagger \mathbf{g}\|_2^2$.

For the sufficiency, $-\theta_1 \geq \sigma \|(H - \theta_1 I_n)^\dagger \mathbf{g}\|_2$ first implies that $-\theta_1 \geq 0$. Assume $\lambda_{\text{opt}} > -\theta_1$.
 Then $\mathbf{h}_{\text{opt}} = -(H + \lambda_{\text{opt}} I_n)^{-1} \mathbf{g}$ and $\mathbf{g} \perp \mathcal{E}_1$ lead to $\mathbf{h}_{\text{opt}} \perp \mathcal{E}_1$, and therefore,

$$\theta_1^2 < \lambda_{\text{opt}}^2 = \sigma^2 \|\mathbf{h}_{\text{opt}}\|_2^2 = \sigma^2 \mathbf{g}^T (H + \lambda_{\text{opt}} I_n)^{-2} \mathbf{g} = \sum_{i=p+1}^n \frac{t_i^2}{(\theta_i + \lambda_{\text{opt}})^2} \leq \sum_{i=p+1}^n \frac{t_i^2}{(\theta_i - \theta_1)^2}.$$

147 But $\sum_{i=p+1}^n \frac{t_i^2}{(\theta_i - \theta_1)^2} = \sigma^2 \|(H - \theta_1 I_n)^{-1} \mathbf{g}\|_2^2$, and by assumption, we have $\theta_1^2 < \lambda_{\text{opt}}^2 \leq \theta_1^2$, a
 148 contradiction. Thus we conclude $\lambda_{\text{opt}} = -\theta_1$. \square

149 By (iii) of Theorem 2.2, as in the trust-region subproblem [9, 1, 20, 41], we define the “hard”
 150 case and “easy” case for the cubic regularization model (2.1).

151 **Definition 2.1.** For (2.1) with $\sigma > 0$, we say it is “hard” case (or the degenerate case) if condition
 152 (2.6) holds. For the hard case, $\lambda_{\text{opt}} = -\theta_1$ and any optimal solution \mathbf{h}_{opt} is of the form

$$153 \quad \mathbf{h}_{\text{opt}} = -(H - \theta_1 I_n)^\dagger \mathbf{g} + \mathbf{u} \quad (2.8)$$

154 where \mathbf{u} is any eigenvector of H corresponding to θ_1 with $\|\mathbf{u}\|_2^2 = \theta_1^2/\sigma^2 - \|(H + \lambda I_n)^\dagger \mathbf{g}\|_2^2$. The
 155 “easy” case (or the non-degenerate case) is characterized by the opposite of the hard case, and
 156 $\lambda_{\text{opt}} > \max(0, -\theta_1)$.

157 3 An associated quadratic eigenvalue problem (QEP)

158 Apart from Newton’s iteration for computing the root λ_{opt} of $q(\lambda) = 0$, analogous to the trust-
 159 region subproblem [1], we can also translate (1.2) into an eigenvalue problem, for which accurate
 160 solution can be efficiently obtained when n is of small- to medium-size. In particular, by (2.4), we
 161 can transform the solution λ_{opt} into the following quadratic eigenvalue problem³ (QEP):

$$162 \quad G(\lambda) := \lambda^2 I_{n+1} + 2\lambda \begin{bmatrix} 0 & 0 \\ 0 & H \end{bmatrix} + \begin{bmatrix} 0 & \sigma \mathbf{g}^\top \\ \sigma \mathbf{g} & H^2 \end{bmatrix} = \begin{bmatrix} \lambda^2 & \sigma \mathbf{g}^\top \\ \sigma \mathbf{g} & (\lambda I_n + H)^2 \end{bmatrix}. \quad (3.1)$$

To see the relation more clearly, using again the spectral decomposition $H = U\Theta U^\top$ and
 denoting $K = \begin{bmatrix} 0 & 1 \\ U & 0 \end{bmatrix}$, we have

$$\begin{aligned} \tilde{G}(\lambda) &:= K^\top G(\lambda) K = \begin{bmatrix} U^\top (\lambda I + H)^2 U & \mathbf{t} \\ \mathbf{t}^\top & \lambda^2 \end{bmatrix} = \begin{bmatrix} (\lambda I + \Theta)^2 & \mathbf{t} \\ \mathbf{t}^\top & \lambda^2 \end{bmatrix} \\ &= \begin{bmatrix} (\lambda + \theta_1)^2 & & & t_1 \\ & \ddots & & \vdots \\ & & (\lambda + \theta_n)^2 & t_n \\ t_1 & \dots & t_n & \lambda^2 \end{bmatrix}. \end{aligned}$$

Denote the eigenvalues of $G(\lambda)$ by $\text{eig}(G(\lambda))$ and $\text{eig}(H) = \{\theta_1, \dots, \theta_n\}$. Noting for all $\lambda \notin \text{eig}(-H)$,
 we have the determinant

$$\begin{aligned} \det G(\lambda) &= \det \tilde{G}(\lambda) = \det \begin{bmatrix} (\lambda I + \Theta)^2 & \mathbf{t} \\ \mathbf{t}^\top & \lambda^2 \end{bmatrix} \\ &= \det(\lambda I + \Theta)^2 \det(\lambda^2 - \mathbf{t}^\top (\lambda I + \Theta)^{-2} \mathbf{t}) \\ &= \left(\lambda^2 - \sum_{i=1}^n \frac{t_i^2}{(\lambda + \theta_i)^2} \right) \prod_{i=1}^n (\lambda + \theta_i)^2 \end{aligned} \quad (3.2)$$

$$= q(\lambda) \prod_{i=1}^n (\lambda + \theta_i)^2 \quad (3.3)$$

$$= \lambda^2 \prod_{i=1}^n (\lambda + \theta_i)^2 - \sum_{i=1}^n t_i^2 \prod_{j \neq i} (\lambda + \theta_j)^2, \quad (3.4)$$

163 and (3.4) is also valid for all λ because (3.4) is a continuous function of λ . This implies that all
 164 the eigenvalues of $G(\lambda)$ which are not in $\text{eig}(-H)$ are the zeros of $q(\lambda) = 0$, i.e., the solutions to
 165 (2.4). This gives the connection between (2.4) and the QEP (3.1).

³A pair (λ, \mathbf{x}) with $\mathbf{x} \neq 0$ is an eigenpair of a polynomial eigenvalue problem $P(\lambda) = P_m \lambda^m + \dots + P_1 \lambda + P_0$ if $P(\lambda)\mathbf{x} = 0$. When $m = 2$ and $m = 1$, we say it is a quadratic and generalized eigenvalue problem, respectively.

3.1 The largest real eigenvalue of $G(\lambda)$ and the associated eigenvector

We next show that the Lagrangian multiplier λ_{opt} can be found by solving the largest real eigenvalue of the QEP (3.1).

Theorem 3.1. For (2.1) with $\sigma > 0$, the Lagrangian multiplier λ_{opt} associated with the global optimal solution \mathbf{h}_{opt} of (2.1) is the largest real eigenvalue of $G(\lambda)$.

Proof. We consider two cases.

Case I: $\mathbf{g} \notin \mathcal{E}_1$ where \mathcal{E}_1 given in (2.6) is the eigenspace associated with the smallest eigenvalue of H . In this case, we know that there exists at least one $t_i \neq 0$ for some $1 \leq i \leq p$, and thus by (3.4), for $\lambda > -\theta_1^+$, we have $\det(G(\lambda)) = q(\lambda) \prod_{i=1}^n (\lambda + \theta_i)^2$ where $q(\lambda)$ is defined in (2.5) and

$$q(\lambda) = \lambda^2 - \sum_{t_i \neq 0, 1 \leq i \leq p} \frac{t_i^2}{(\theta_1 + \lambda)^2} - \sum_{i=p+1}^n \frac{t_i^2}{(\theta_i + \lambda)^2} \begin{cases} \rightarrow +\infty, & \text{as } \lambda \rightarrow +\infty \\ < 0 & \text{as } \lambda \rightarrow 0 \\ \rightarrow -\infty, & \text{as } \lambda \rightarrow -\theta_1 \end{cases}. \quad (3.5)$$

Therefore, together with the monotonicity of $q(\lambda)$ on $(-\theta_1^+, +\infty)$, we know that there is a unique solution in $(-\theta_1^+, +\infty)$ for $\det(G(\lambda)) = 0$, which by Theorem 2.2 is λ_{opt} .

Case II: $\mathbf{g} \perp \mathcal{E}_1$. In this case, by (3.4), we know that $-\theta_1 \in \text{eig}(G(\lambda))$, and also $t_i = 0$ for all $1 \leq i \leq p$. If $\mathbf{g} = \mathbf{0}$, then $\text{eig}(G(\lambda)) = \{-\theta_1, -\theta_1, \dots, -\theta_n, -\theta_n, 0\}$, and the largest eigenvalue is $-\theta_1^+$. According to Theorem 2.2 (i), the claim of this theorem is true.

For $\mathbf{g} \neq \mathbf{0}$, by (3.4), for $\lambda > -\theta_1^+$, we have $\det(G(\lambda)) = \tilde{q}(\lambda) \prod_{i=1}^n (\lambda + \theta_i)^2$ where $\tilde{q}(\lambda)$ is given by (2.7). Note that $\tilde{q}(\lambda) \rightarrow +\infty$ as $\lambda \rightarrow +\infty$, and is monotonically increasing on $[-\theta_1^+, +\infty)$. If $-\theta_1 \leq 0$, then $-\theta_1^+ = 0$ and by noting $\tilde{q}(0) < 0$, we know that there is a unique solution in $(-\theta_1^+, +\infty)$ for $\det(G(\lambda)) = 0$, which by Theorem 2.2 is λ_{opt} . Otherwise, $-\theta_1 > 0$ and $-\theta_1^+ = -\theta_1$. Note that

$$\tilde{q}(-\theta_1) = \theta_1^2 - \sum_{i=p+1}^n \frac{t_i^2}{(\theta_i - \theta_1)^2} = \theta_1^2 - \sigma^2 \|(H - \theta_1 I_n)^\dagger \mathbf{g}\|_2^2.$$

Therefore, if $\tilde{q}(-\theta_1) < 0$, then there is still a unique solution in $(-\theta_1^+, +\infty)$ for $\det(G(\lambda)) = 0$, which by Theorem 2.2 is just λ_{opt} . But if $\tilde{q}(-\theta_1) \geq 0$, then there is no solution in $(-\theta_1^+, +\infty)$ for $\det(G(\lambda)) = 0$, and the largest real eigenvalue of $G(\lambda)$ is therefore $-\theta_1$. The latter case corresponds to the hard case. \square

We next show that the solution \mathbf{h}_{opt} of (2.1) can be obtained from the eigenpair $(\lambda_{\text{opt}}, \mathbf{z})$ associated with the largest real eigenvalue λ_{opt} of $G(\lambda)$.

Theorem 3.2. For (2.1) with $\sigma > 0$, suppose $\mathbf{z} = [\alpha; \mathbf{y}] \in \mathbb{R}^{n+1}$ is the normalized eigenvector of $G(\lambda)$ associated with the largest real eigenvalue λ_{opt} . Then we have

(i) if $\alpha \neq 0$, then $\mathbf{h}_{\text{opt}} = (H\mathbf{y} + \lambda_{\text{opt}}\mathbf{y})/(\sigma\alpha)$ is the solution to (2.1);

(ii) if $\alpha = 0$, then $\lambda_{\text{opt}} = -\theta_1$, \mathbf{y} is an eigenvector of H associated with $-\lambda_{\text{opt}}$ and $\mathbf{h}_{\text{opt}} = (H + \lambda_{\text{opt}}I_n)^\dagger \mathbf{g} + \eta\mathbf{y}$ is the solution to (2.1), where

$$\eta = \pm \sqrt{\lambda_{\text{opt}}^2 - \sigma^2 \|(H + \lambda_{\text{opt}}I_n)^\dagger \mathbf{g}\|_2^2}.$$

Proof. By $G(\lambda_{\text{opt}})\mathbf{z} = \mathbf{0}$, we have

$$\lambda_{\text{opt}}^2 \alpha + \sigma \mathbf{g}^T \mathbf{y} = 0, \quad \text{and} \quad \sigma \alpha \mathbf{g} + (H + \lambda_{\text{opt}}I_n)^2 \mathbf{y} = \mathbf{0}. \quad (3.6)$$

For (i), we first consider the case $\lambda_{\text{opt}} > -\theta_1$, which by the second equation in (3.6) implies $\mathbf{h}_{\text{opt}} = (H\mathbf{y} + \lambda_{\text{opt}}\mathbf{y})/(\sigma\alpha)$ satisfies $(H + \lambda_{\text{opt}}I_n)\mathbf{h}_{\text{opt}} = -\mathbf{g}$; moreover, by (3.6) again, we know that $\sigma \|\mathbf{h}_{\text{opt}}\|_2 = \lambda$, and according to Theorem 2.1, the claim is true. For the case $\lambda_{\text{opt}} = -\theta_1$, Theorem 2.2 (iii) indicates that $\mathbf{g} \perp \mathcal{E}_1$, and thus, the second equation in (3.6) implies

$$-(H + \lambda_{\text{opt}}I_n)^\dagger \mathbf{g} = (H\mathbf{y} + \lambda_{\text{opt}}\mathbf{y})/(\sigma\alpha).$$

192 Moreover, by (3.6) and $\alpha \neq 0$, we have $\|(H + \lambda_{\text{opt}}I_n)^\dagger \mathbf{g}\|_2 = -\theta_1/\sigma$, and hence $\mathbf{h}_{\text{opt}} = (H\mathbf{y} +$
193 $\lambda_{\text{opt}}\mathbf{y})/(\sigma\alpha)$ is the solution to (2.1). We remark that this case ($\lambda_{\text{opt}} = -\theta_1$ and $\alpha \neq 0$) only
194 happens if $\mathbf{g} \perp \mathcal{E}_1$ and $-\theta_1 = \sigma\|(H - \theta_1I_n)^\dagger \mathbf{g}\|_2$.

195 For (ii), we first know $\|\mathbf{y}\|_2 = 1$, and also (3.6) reduces to $\sigma\mathbf{g}^\top \mathbf{y} = 0$ and $(H + \lambda I_n)^2 \mathbf{g} = \mathbf{0}$. The
196 largest real eigenvalue of $G(\lambda)$ in this case must be $-\theta_1$ because otherwise $\mathbf{y} = \mathbf{0}$. Thus, $\mathbf{y} \in \mathcal{E}_1$
197 and by Theorem 2.2 (iii), we know that the solution in this case is $\mathbf{h}_{\text{opt}} = -(H + \lambda_{\text{opt}}I_n)^\dagger \mathbf{g} + \eta \mathbf{y}$
198 with η as claimed. \square

In practice, similar to [1], when $|\alpha| \neq 0$ is close to zero, we can use

$$\mathbf{h}_{\text{opt}} = \text{sign}(\alpha) \left| \frac{\lambda_{\text{opt}}}{\sigma} \right| \frac{H\mathbf{y} + \lambda_{\text{opt}}\mathbf{y}}{\|H\mathbf{y} + \lambda_{\text{opt}}\mathbf{y}\|_2}$$

199 instead of $\mathbf{h}_{\text{opt}} = (H\mathbf{y} + \lambda_{\text{opt}}\mathbf{y})/(\sigma\alpha)$ to compute the solution \mathbf{h}_{opt} .

200 3.2 All the eigenvalues of $G(\lambda)$

201 We take a close look at the eigenvalues of $G(\lambda)$ in this subsection. We separate our discussion into
202 two scenarios: the generic case when θ_i are distinct and the remaining special case.

203 3.2.1 The generic case: θ_i are distinct

In this situation, by (3.4), we have $\det(G(-\theta_i)) = -t_i^2 \prod_{j \neq i} (\theta_j - \theta_i)^2$, and it is clear that

$$t_i = 0 \iff \det(G(-\theta_i)) = 0 \iff -\theta_i \in \text{eig}(G(\lambda)).$$

When all $t_i \neq 0$ for all $i = 1, 2, \dots, n$, then $\text{eig}(G(\lambda)) \cap \text{eig}(-H) = \emptyset$, and moreover, $0 \notin$
 $\text{eig}(G(\lambda))$, because, otherwise by (3.4), $\det(G(0)) = -\sum_{i=1}^n t_i^2 \prod_{i=1}^n \theta_i^2 \neq 0$, implying $0 \notin \text{eig}(G(\lambda))$,
a contradiction. Write $q(\lambda) = \lambda^2 \delta(\lambda)$, where

$$\delta(\lambda) := 1 - \sum_{i=1}^n \frac{t_i^2}{\lambda^2(\lambda + \theta_i)^2}, \quad \text{and} \quad \delta'(\lambda) = 2 \sum_{i=1}^n \frac{(2\lambda + \theta_i)t_i^2}{\lambda^3(\lambda + \theta_i)^3}.$$

204 According to (3.3), the zeros of $\delta(\lambda)$ are the eigenvalues of $G(\lambda)$. To have a clear picture of the
205 distribution of the eigenvalues, based on the signs of $\delta(\lambda)$ and $\delta'(\lambda)$, we draw a diagram of the func-
206 tion $\delta(\lambda)$ in Figure 3.1 as an illustration for the case $\theta_1 > 0$. In this illustration, the branches in the
207 intervals $(-\theta_n, -\theta_{n-1})$, $(-\theta_{n-1}, -\theta_{n-2})$, $(-\theta_1, 0)$ represent three types of eigenvalues respectively:
208 two different real eigenvalues, two same real eigenvalues, and two conjugate complex eigenvalues.
209 Together with $(-\infty, -\theta_n)$ and $(0, \infty)$, these $n + 2$ intervals contain $2n + 2$ zeros in total. Clearly
210 there exists unique real eigenvalue of $G(\lambda)$ larger than $\max(0, -\theta_1)$, and according to Theorem 2.1,
211 this real eigenvalue is just λ_{opt} . We remark that this case, i.e., $t_i \neq 0$ for all i , is the easy case of
212 (2.1) by Definition 2.1.

213 For the case when $t_i = 0$ for $i \in \mathcal{I} \subset \{1, \dots, n\}$, we know that $-\theta_i$ is an eigenvalue of $G(\lambda)$.
214 To obtain other eigenvalues, we consider $\tilde{G}_{\mathcal{I}^c}(\lambda) = I_{\mathcal{I}^c}^\top \tilde{G}(\lambda) I_{\mathcal{I}^c}$, where $I_{\mathcal{I}^c}$ is a matrix consisting of
215 \mathbf{e}_j for $j \in \{1, \dots, n\} \setminus \mathcal{I}$. The eigenvalues of $\tilde{G}_{\mathcal{I}^c}(\lambda)$ can be treated in the case above. If $|\mathcal{I}| = n$,
216 namely $\mathbf{g} = \mathbf{t} = \mathbf{0}$ (implying that 0 is also an eigenvalue of $G(\lambda)$), there is no real eigenvalue of
217 $G(\lambda)$ larger than $\max(0, -\theta_1)$, and according to Theorem 2.1, $\lambda_{\text{opt}} = \max(0, -\theta_1)$. Otherwise,
218 $|\mathcal{I}| < n$, and there exists a unique real eigenvalue, say χ , of $G(\lambda)$ larger than $\max(0, -\theta_i)_{i \notin \mathcal{I}}$. For
219 this case, if $1 \notin \mathcal{I}$, then χ is the largest real eigenvalue of $G(\lambda)$ which is just λ_{opt} (the easy case)
220 according to Theorem 2.1, while $1 \in \mathcal{I}$, χ and $-\theta_1$ are two eigenvalues. The latter case leads to
221 the easy case when $\chi > -\theta_1$ and the hard case $\chi \leq -\theta_1$ by Definition 2.1.

222 3.2.2 The special case

223 We now assume that the distinct values of θ_i are μ_1, \dots, μ_m with $m < n$. First we consider the
224 subset $\mathcal{I}_1 \subset \{1, \dots, n\}$ where $\theta_i = \mu_1$ for any $i \in \mathcal{I}_1$.

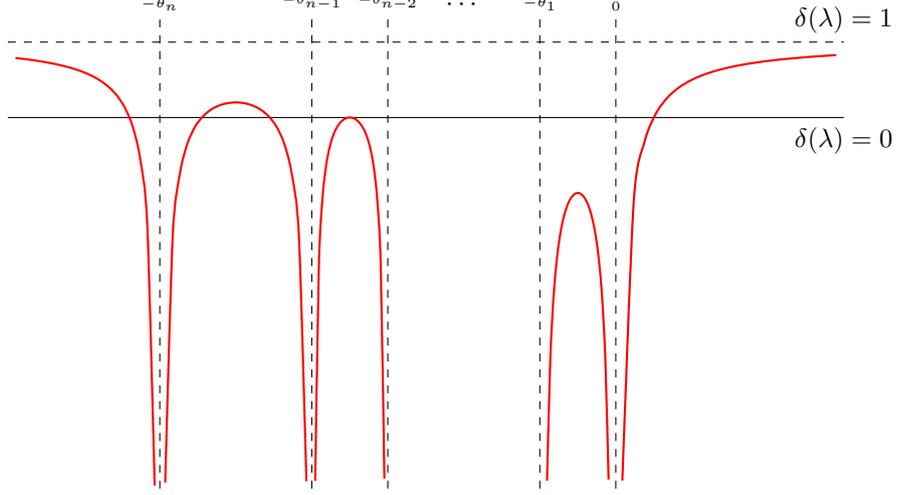


Figure 3.1: $\delta(\lambda)$ on the case $\theta_1 > 0$

225 When $t_i = 0$ for all $i \in \mathcal{I}_1$, then by (3.2)

$$226 \quad \det(G(\lambda)) = \left(\lambda^2 - \sum_{i \notin \mathcal{I}_1} \frac{t_i^2}{(\lambda + \theta_i)^2} \right) (\lambda + \mu_1)^{2|\mathcal{I}_1|} \prod_{i \notin \mathcal{I}_1} (\lambda + \theta_i)^2.$$

227 Thus, $-\mu_1$ is an eigenvalue of $G(\lambda)$ with (algebraic) multiplicity $2|\mathcal{I}_1|$. Analogous to the generic
 228 case, the other eigenvalues can be obtained by considering $\tilde{G}_{\mathcal{I}_1^c}(\lambda) = I_{\mathcal{I}_1^c}^T \tilde{G}(\lambda) I_{\mathcal{I}_1^c}$.

229 Otherwise (i.e., $\sum_{i \in \mathcal{I}_1} t_i^2 \neq 0$), we have by (3.2) that

$$230 \quad \det(G(\lambda)) = \left(\lambda^2 (\lambda + \mu_1)^2 - \sum_{i \notin \mathcal{I}_1} \frac{(\lambda + \mu_1)^2 t_i^2}{(\lambda + \theta_i)^2} - \sum_{i \in \mathcal{I}_1} t_i^2 \right) (\lambda + \mu_1)^{2(|\mathcal{I}_1| - 1)} \prod_{i \notin \mathcal{I}_1} (\lambda + \theta_i)^2.$$

231 Thus, $-\mu_1$ is an eigenvalue of $G(\lambda)$ with (algebraic) multiplicity $2(|\mathcal{I}_1| - 1)$, and the other eigen-
 232 values can be obtained by considering $\tilde{G}_{\mathcal{I}_1^c}(\lambda) = \begin{bmatrix} (\lambda + \mu_1)^2 & s_1 e_n^T \\ s_1 e_n & I_{\mathcal{I}_1^c}^T \tilde{G}(\lambda) I_{\mathcal{I}_1^c} \end{bmatrix}$, where $s_1^2 = \sum_{i \in \mathcal{I}_1} t_i^2$.

233 The above arguments can be continuously applied to μ_2, \dots, μ_m to obtain all the eigenvalues
 234 of $G(\lambda)$ and the details are omitted.

235 3.3 Associated generalized eigenvalue problems

236 In a recent work by Lieder [25], it is shown that the optimal λ_{opt} is the largest real eigenvalue of
 237 the following generalized eigenvalue problem:

$$238 \quad \mathcal{M}(\lambda) = \begin{bmatrix} 0 & 0 & 0 & -\mathbf{g}^T \\ 0 & \sigma I_n & 0 & -H \\ 0 & 0 & \sigma & 0 \\ -\mathbf{g} & -H & 0 & 0 \end{bmatrix} - \lambda \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & I_n \\ 1 & 0 & 0 & 0 \\ 0 & I_n & 0 & 0 \end{bmatrix} \in \mathbb{R}^{2(n+1) \times 2(n+1)}. \quad (3.7)$$

239 It is known that a QEP can be linearized to various types of generalized eigenvalue problem (see
 240 e.g., [36]). Thus, our QEP (3.1) can lead us to other generalized eigenvalue problems. For example,

$$241 \quad \mathcal{L}(\lambda) = \begin{bmatrix} -A & 0 \\ 0 & I_{n+1} \end{bmatrix} - \lambda \begin{bmatrix} B & I_{n+1} \\ I_{n+1} & 0 \end{bmatrix} =: \mathcal{A} - \lambda \mathcal{B} \in \mathbb{R}^{2(n+1) \times 2(n+1)} \quad (3.8)$$

where

$$B = \begin{bmatrix} 0 & 0 \\ 0 & 2H \end{bmatrix}, \quad A = \begin{bmatrix} 0 & \sigma \mathbf{g}^T \\ \sigma \mathbf{g} & H^2 \end{bmatrix}.$$

242 We will see in the next theorem that the eigenpair of $G(\lambda)$ associated with λ_{opt} can be equivalently
243 obtained via a generalized eigenvalue problem.

244 **Theorem 3.3.** ([24, Theorem 4.1]) *Let $\mathcal{L}(\lambda)$ be given by (3.8), then we have*

- 245 (1) *The set of eigenvalues of $G(\lambda)$ is the same as that of the matrix pencil $A - \lambda B$.*
246 (2) *The inertia (i.e., the number of the positive, zero, and negative eigenvalues respectively) of*
247 *\mathcal{B} is $(n + 1, 0, n + 1)$.*
248 (3) *If (λ, \mathbf{z}) is an eigenpair of $G(\lambda)$, then $\left(\lambda, \begin{bmatrix} \mathbf{z} \\ \lambda \mathbf{z} \end{bmatrix}\right)$ is an eigenpair of $\mathcal{L}(\lambda)$.*
249 (4) *If $\left(\lambda, \begin{bmatrix} \mathbf{z} \\ \mathbf{t} \end{bmatrix}\right)$ is an eigenpair of $\mathcal{L}(\lambda)$, then (λ, \mathbf{z}) is an eigenpair of $G(\lambda)$ and $\mathbf{t} = \lambda \mathbf{z}$.*

250 **Remark 3.1.** It should be pointed out that even though both (3.7) and (3.8) are $2(n + 1)$ dimen-
251 sional generalized eigenvalue problems, (3.7) is preferable as it does not involve H^2 . However, our
252 QEP (3.1) has the following advantages.

- 253 1) The QEP (3.1) is of dimension $(n + 1)$ and efficient Krylov subspace methods working on
254 \mathbb{R}^{n+1} such as the second-order Arnoldi process (SOAR) [2] can be directly applied to solve
255 (1.2).
256 2) The QEP (3.1) is **more flexible**. First, there are various types of generalized eigenvalue
257 problems that can be derived from (3.1) by linearization [36]. For instance, by taking the
258 advantage of the coefficient matrix I_{n+1} in the term λ^2 in (3.1), another commonly used
259 linearization in the literature [36] leads to the following standard eigenvalue problem:

260
$$C\mathbf{y} = \lambda\mathbf{y}, \quad C = \begin{bmatrix} -B & -A \\ I_{n+1} & 0 \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} \lambda \mathbf{x} \\ \mathbf{x} \end{bmatrix}.$$

Secondly, in many applications of the cubic regularization of Newton method for (1.1), the
BFGS update [28] will be used to approximate the inverse of the Hessian matrix, where only
 $\hat{H} = H^{-1}$ is available. In this situation, by noting

$$\det \left(\underbrace{\begin{bmatrix} 1 & & & \\ & H^{-1} & & \\ & & G(\lambda) & \\ & & & 1 & & \\ & & & & H^{-1} & \end{bmatrix}}_{=:\hat{G}(\lambda)} \right) = \det(H^{-2}) \cdot \det(G(\lambda)),$$

261 we have a QEP involving just \hat{H} :

262
$$\hat{G}(\lambda) = \lambda^2 \begin{bmatrix} 0 & 0 \\ 0 & \hat{H}^2 \end{bmatrix} + 2\lambda \begin{bmatrix} 0 & 0 \\ 0 & \hat{H} \end{bmatrix} + \begin{bmatrix} 0 & \sigma \mathbf{g}^T \hat{H} \\ \sigma \hat{H} \mathbf{g} & I_n \end{bmatrix},$$

263 and the optimal λ_{opt} can be computed by finding the largest real eigenvalue of $\hat{G}(\lambda)$.

264 4 A Lanczos method

265 We next discuss a Lanczos type procedure introduced in [7, Section 6.2] for solving (1.2). The
266 approach is analogous to the generalized Lanczos trust-region (GLTR) method proposed in [15] for
267 the trust-region problem. GLTR is an efficient Lanczos type method for large-scale minimization

268 problems and its convergence analysis was recently established in [41, 42] and efficient restarting
 269 techniques are developed in [40]. For the cubic regularization model (2.1), the approach begins
 270 with forming an ℓ -th Krylov subspace $\mathcal{K}_\ell(H, \mathbf{g}) = \text{span}(Q_\ell)$ via the standard Lanczos process
 271 (Algorithm 1) for ℓ less than the grade τ of \mathbf{g} with respect to H (i.e., τ is the smallest number
 272 that the Lanczos process breaks at step 6), and we have

$$273 \quad HQ_\ell = Q_\ell S_\ell + \gamma_\ell \mathbf{q}_{\ell+1} \mathbf{e}_\ell^\top, \quad Q_\ell \mathbf{e}_1 = \mathbf{g} / \|\mathbf{g}\|_2, \quad (4.1)$$

274 where

$$275 \quad S_\ell = Q_\ell^\top H Q_\ell = \begin{bmatrix} \delta_0 & \gamma_1 & & & \\ \gamma_1 & \delta_1 & \ddots & & \\ & \ddots & \ddots & \ddots & \\ & & & \gamma_\ell & \delta_\ell \end{bmatrix}$$

is tridiagonal, $Q_\ell = [\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_\ell]$ is orthogonal.

Algorithm 1 The standard Lanczos Process

It computes an orthogonal basis matrix Q_ℓ of $\mathcal{K}_\ell(H, \mathbf{g})$.

```

1: set  $\mathbf{q}_0 = 0, \gamma_0 = 0, \mathbf{q}_1 = \mathbf{g} / \|\mathbf{g}\|_2$ ;
2: for  $j = 1, 2, \dots, \ell$  do
3:    $\mathbf{t} = H\mathbf{q}_j, \delta_{j-1} = \mathbf{q}_j^\top \mathbf{t}$ ;
4:    $\mathbf{t} = \mathbf{t} - \delta_{j-1} \mathbf{q}_j - \gamma_{j-1} \mathbf{q}_{j-1}, \gamma_j = \|\mathbf{t}\|_2$ ;
5:   if  $\gamma_j = 0$  then
6:     break;
7:   else
8:      $\mathbf{q}_{j+1} = \mathbf{t} / \gamma_j$ ;
9:   end if
10: end for

```

276 The approximation $\mathbf{h}_\ell \in \mathbb{R}^n$ at the ℓ -th step is obtained by

$$277 \quad \mathbf{h}_\ell = \underset{\mathbf{h} \in \mathcal{K}_\ell(H, \mathbf{g})}{\operatorname{argmin}} m(\mathbf{h}), \quad (4.2)$$

279 which, by denoting $\mathbf{h} \in \mathcal{K}_\ell(H, \mathbf{g})$ by $\mathbf{h} = Q_\ell \mathbf{s}$ with $\mathbf{s} \in \mathbb{R}^\ell$, can be solved equivalently as the
 280 following smaller sized cubic regularization model:

$$281 \quad \min_{\mathbf{s} \in \mathbb{R}^\ell} \left\{ \|\mathbf{g}\|_2 \mathbf{e}_1^\top \mathbf{s} + \frac{1}{2} \mathbf{s}^\top S_\ell \mathbf{s} + \frac{\sigma}{3} \|\mathbf{s}\|_2^3 \right\}. \quad (4.3)$$

282 Denote the solution of (4.3) by $\mathbf{s}_{\text{opt}, \ell}$, which can be obtained by solving the associated eigenvalue
 283 problem discussed in Theorem 3.2. Also, we have $\mathbf{h}_\ell = Q_\ell \mathbf{s}_{\text{opt}, \ell}$, and the corresponding Lagrangian
 284 multiplier is

$$285 \quad \lambda_\ell = \sigma \|\mathbf{h}_\ell\|_2 = \sigma \|\mathbf{s}_{\text{opt}, \ell}\|_2. \quad (4.4)$$

286 4.1 Convergence

287 To establish the convergence results of $\{m(\mathbf{h}_\ell)\}$ and $\{\mathbf{h}_\ell\}$, we first have the following lemmas.

288 **Lemma 4.1.** *Suppose $\mathbf{g} \notin \mathcal{E}_1$ where \mathcal{E}_1 is the eigenspace of H associated with the smallest eigen-*
 289 *value θ_1 . Then we have*

- 290 (i) *Before the breakdown at step 6 in Algorithm 1, it holds that $S_\ell + \lambda_\ell I_\ell \succ 0$, where S_ℓ is the*
 291 *tridiagonal matrix given in (4.1) and λ_ℓ is given by (4.4);*
- 292 (ii) *If the standard Lanczos process (Algorithm 1) breaks at the τ -th step, then $\lambda_{\text{opt}} = \lambda_\tau$ and*
 293 *$\mathbf{h}_\tau = \mathbf{h}_{\text{opt}}$.*

294 *Proof.* For i), we note that S_ℓ is irreducible, and then by [15, Theorem 5.3] (see also [29, Theorem
 295 7.9.5]), any eigenvector \mathbf{v} of S_ℓ satisfies $\mathbf{v}^\top \mathbf{e}_1 \neq 0$. Thus, by Theorem 2.2 (iii), we know that (4.3)
 296 is an easy case and the associated λ_ℓ is strictly larger than the smallest eigenvalue of S_ℓ .

We prove ii) by showing that $\lambda_{\text{opt}} = \lambda_\tau$. According to Theorem 2.1, the assumption of $\mathbf{g} \notin \mathcal{E}_1$
 and the arguments in (3.5), it follows that λ_{opt} and λ_τ are the largest real roots to the systems

$$q(\lambda) = \lambda^2 - \sigma^2 \|(H + \lambda I_n)^{-1} \mathbf{g}\|_2^2 = 0, \quad \text{and} \quad \widehat{q}(\lambda) = \lambda^2 - \sigma^2 \|\mathbf{g}\|_2 \|(S_\tau + \lambda I_\tau)^{-1} \mathbf{e}_1\|_2^2 = 0,$$

respectively. Indeed, in this case we can show that $q(\lambda) = \widehat{q}(\lambda)$ for any $\lambda \notin \text{eig}(H)$, and hence
 $\lambda_{\text{opt}} = \widehat{\lambda}_\tau$. To this end, by assumptions, we know that $\text{span}(Q_\tau) = \mathcal{K}_\tau(H, \mathbf{g})$ is an invariant
 subspace of H (implying $\text{eig}(S_\tau) \subseteq \text{eig}(H)$), which contains the eigenvectors associated with the
 smallest eigenvalue θ_1 of H . So the smallest eigenvalue of S_τ is θ_1 and $HQ_\tau = Q_\tau S_\tau$ leading to

$$(H + \lambda I_n)Q_\tau = Q_\tau(S_\tau + \lambda I_\tau).$$

Thus for any $\lambda \notin \text{eig}(H)$, we get from $Q_\tau \mathbf{e}_1 = \mathbf{g}/\|\mathbf{g}\|_2$ that

$$(H + \lambda I_n)^{-1} \mathbf{g} = (H + \lambda I_n)^{-1} Q_\tau \mathbf{e}_1 \|\mathbf{g}\|_2 = Q_\tau (S_\tau + \lambda I_\tau)^{-1} \mathbf{e}_1 \|\mathbf{g}\|_2$$

297 which leads to $q(\lambda) = \widehat{q}(\lambda)$, and the conclusion follows. \square

298 The linear convergence of λ_ℓ in (4.4) to λ_{opt} has been previously discussed in [5, 6], in which
 299 the proof follows similarly as that for the trust-region subproblem. Here, following the argument in
 300 [41] for the trust-region subproblem, we provide a different way which can render sharper bounds
 301 (refer to Remark 4.2) for the approximate objective function value as well as the solution.

Lemma 4.2. *Under the assumption of Lemma 4.1, let \mathcal{L} be any subspace of \mathbb{R}^n with $\dim(\mathcal{L}) \geq 1$
 and \mathbf{c} be the solution to $\min_{\mathbf{h} \in \mathcal{L}} m(\mathbf{h})$. Then for any nonzero vector $\mathbf{h} \in \mathcal{L}$, we have*

$$0 \leq m(\mathbf{c}) - m(\mathbf{h}_{\text{opt}}) \leq 2\|H_{\text{opt}}\|_2 \|\mathbf{h} - \mathbf{h}_{\text{opt}}\|_2^2, \quad (4.5)$$

$$\|\mathbf{c} - \mathbf{h}_{\text{opt}}\|_2 \leq 2\sqrt{\kappa} \|\mathbf{h} - \mathbf{h}_{\text{opt}}\|_2, \quad (4.6)$$

302 where $H_{\text{opt}} = H + \lambda_{\text{opt}} I_n$ and $\kappa = \frac{\theta_n + \lambda_{\text{opt}}}{\theta_1 + \lambda_{\text{opt}}}$.

Proof. Let $\mathbf{h} \in \mathcal{L}$ be any nonzero vector and let $\mathbf{v} = \mathbf{h} \frac{\lambda_{\text{opt}}}{\|\mathbf{h}\|_2 \sigma} \in \mathcal{L}$. Define $\mathbf{m} = \mathbf{v} - \mathbf{h}_{\text{opt}}$. Since \mathbf{c} is
 the minimizer, we have

$$\begin{aligned} 0 &\leq m(\mathbf{c}) - m(\mathbf{h}_{\text{opt}}) \leq m(\mathbf{v}) - m(\mathbf{h}_{\text{opt}}) \\ &= \mathbf{g}^\top \mathbf{m} + \mathbf{m}^\top H \mathbf{h}_{\text{opt}} + \frac{1}{2} \mathbf{m}^\top H \mathbf{m} + \frac{\sigma}{3} (\|\mathbf{v}\|_2^3 - \|\mathbf{h}_{\text{opt}}\|_2^3) \\ &= -\lambda_{\text{opt}} \mathbf{m}^\top \mathbf{h}_{\text{opt}} + \frac{1}{2} \mathbf{m}^\top H \mathbf{m}, \quad (H \mathbf{h}_{\text{opt}} = -\mathbf{g} - \lambda_{\text{opt}} \mathbf{h}_{\text{opt}}, \|\mathbf{v}\|_2 = \|\mathbf{h}_{\text{opt}}\|_2) \\ &= \frac{1}{2} \mathbf{m}^\top (H + \lambda_{\text{opt}} I_n) \mathbf{m} \\ &\leq \frac{1}{2} \|H_{\text{opt}}\|_2 \|\mathbf{m}\|_2^2, \end{aligned} \quad (4.7)$$

where the last equality is due to the fact $\mathbf{m}^\top \mathbf{h}_{\text{opt}} = -\|\mathbf{m}\|_2^2/2$ which follows from

$$\mathbf{h}_{\text{opt}}^\top \mathbf{h}_{\text{opt}} = \mathbf{v}^\top \mathbf{v} = (\mathbf{h}_{\text{opt}} + \mathbf{m})^\top (\mathbf{h}_{\text{opt}} + \mathbf{m}) = \mathbf{h}_{\text{opt}}^\top \mathbf{h}_{\text{opt}} + 2\mathbf{m}^\top \mathbf{h}_{\text{opt}} + \|\mathbf{m}\|_2^2.$$

Furthermore, $\|\mathbf{m}\|_2 = \|\mathbf{v} - \mathbf{h}_{\text{opt}}\|_2 \leq \|\mathbf{v} - \mathbf{h}\|_2 + \|\mathbf{h} - \mathbf{h}_{\text{opt}}\|_2$, and

$$\begin{aligned} \|\mathbf{v} - \mathbf{h}\|_2 &= \left\| \mathbf{h} - \mathbf{h} \frac{\lambda_{\text{opt}}}{\|\mathbf{h}\|_2 \sigma} \right\|_2 \\ &= \|\mathbf{h}\|_2 \cdot \left\| 1 - \frac{\lambda_{\text{opt}}}{\|\mathbf{h}\|_2 \sigma} \right\|_2 \end{aligned}$$

$$\begin{aligned}
&= \left\| \|\mathbf{h}\|_2 - \frac{\lambda_{\text{opt}}}{\sigma} \right\|_2 \\
&= \left\| \|\mathbf{h}\|_2 - \|\mathbf{h}_{\text{opt}}\|_2 \right\|_2, \quad (\|\mathbf{h}_{\text{opt}}\|_2 = \frac{\lambda_{\text{opt}}}{\sigma}) \\
&\leq \|\mathbf{h} - \mathbf{h}_{\text{opt}}\|_2.
\end{aligned}$$

303 Thus, $\|\mathbf{m}\|_2 \leq 2\|\mathbf{h} - \mathbf{h}_{\text{opt}}\|_2$, and together with (4.7), the claim (4.5) follows.

304 Now we consider (4.6). Denote by $Q_{\mathcal{L}}$ an orthonormal basis of \mathcal{L} . Let

$$305 \quad \mathbf{c} = Q_{\mathcal{L}}\tilde{\mathbf{h}}_{\mathcal{L}} = \underset{\mathbf{h} \in \mathcal{L}}{\text{argmin}} m(\mathbf{h}), \quad \text{where } \tilde{\mathbf{h}}_{\mathcal{L}} = \underset{\tilde{\mathbf{h}} \in \mathbb{R}^{\dim(\mathcal{L})}}{\text{argmin}} m(Q_{\mathcal{L}}\tilde{\mathbf{h}}).$$

306 By Theorem 2.1, we know

$$307 \quad (Q_{\mathcal{L}}^T H Q_{\mathcal{L}} + \lambda_{\mathcal{L}} I)\tilde{\mathbf{h}}_{\mathcal{L}} = -Q_{\mathcal{L}}^T \mathbf{g}, \quad \lambda_{\mathcal{L}} = \sigma \|\tilde{\mathbf{h}}_{\mathcal{L}}\|_2, \quad Q_{\mathcal{L}}^T H Q_{\mathcal{L}} + \lambda_{\mathcal{L}} I \succeq 0,$$

where $\lambda_{\mathcal{L}}$ is the corresponding Lagrangian multiplier satisfying $\lambda_{\mathcal{L}} = \sigma \|\mathbf{c}\|_2$. Since

$$\begin{aligned}
&m(\mathbf{c}) - m(\mathbf{h}_{\text{opt}}) \\
&= \mathbf{g}^T(\mathbf{c} - \mathbf{h}_{\text{opt}}) + \frac{1}{2}\mathbf{c}^T H \mathbf{c} - \frac{1}{2}\mathbf{h}_{\text{opt}}^T H \mathbf{h}_{\text{opt}} + \frac{\sigma}{3}\|\mathbf{c}\|_2^3 - \frac{\sigma}{3}\|\mathbf{h}_{\text{opt}}\|_2^3 \\
&= -\mathbf{h}_{\text{opt}}^T (H + \lambda_{\text{opt}} I)(\mathbf{c} - \mathbf{h}_{\text{opt}}) + \frac{1}{2}\mathbf{c}^T (H + \lambda_{\text{opt}} I)\mathbf{c} - \frac{1}{2}\mathbf{h}_{\text{opt}}^T (H + \lambda_{\text{opt}} I)\mathbf{h}_{\text{opt}} \\
&\quad - \frac{1}{2}\lambda_{\text{opt}}\|\mathbf{c}\|_2^2 + \frac{1}{2}\lambda_{\text{opt}}\|\mathbf{h}_{\text{opt}}\|_2^2 + \frac{\sigma}{3}\|\mathbf{c}\|_2^3 - \frac{\sigma}{3}\|\mathbf{h}_{\text{opt}}\|_2^3 \\
&= \frac{1}{2}(\mathbf{c} - \mathbf{h}_{\text{opt}})^T (H + \lambda_{\text{opt}} I)(\mathbf{c} - \mathbf{h}_{\text{opt}}) \\
&\quad - \frac{1}{2}\lambda_{\text{opt}} \left(\frac{\lambda_{\mathcal{L}}}{\sigma} \right)^2 + \frac{1}{2}\lambda_{\text{opt}} \left(\frac{\lambda_{\text{opt}}}{\sigma} \right)^2 + \frac{\sigma}{3} \left(\frac{\lambda_{\mathcal{L}}}{\sigma} \right)^3 - \frac{\sigma}{3} \left(\frac{\lambda_{\text{opt}}}{\sigma} \right)^3 \\
&= \frac{1}{2}(\mathbf{c} - \mathbf{h}_{\text{opt}})^T (H + \lambda_{\text{opt}} I)(\mathbf{c} - \mathbf{h}_{\text{opt}}) + \frac{1}{6\sigma^2}(\lambda_{\text{opt}}^3 - 3\lambda_{\text{opt}}\lambda_{\mathcal{L}}^2 + 2\lambda_{\mathcal{L}}^3) \\
&= \frac{1}{2}(\mathbf{c} - \mathbf{h}_{\text{opt}})^T (H + \lambda_{\text{opt}} I)(\mathbf{c} - \mathbf{h}_{\text{opt}}) + \frac{1}{6\sigma^2}(\lambda_{\text{opt}} - \lambda_{\mathcal{L}})^2(\lambda_{\text{opt}} + 2\lambda_{\mathcal{L}}) \\
&\geq \frac{1}{2}(\theta_1 + \lambda_{\text{opt}})\|\mathbf{c} - \mathbf{h}_{\text{opt}}\|_2^2,
\end{aligned}$$

308 together with (4.5), we have (4.6). □

Lemma 4.3 (Bernstein [3]). *Given $\phi > 1$, the best approximating polynomial $p_{\ell}(x) \in \mathbb{P}_{\ell}$ of $\frac{1}{x-\phi}$ in $[-1, 1]$ satisfies*

$$\frac{1}{x-\phi} - p_{\ell}(x) = \frac{(\phi + \sqrt{\phi^2 - 1})^{-\ell}}{\phi^2 - 1} \cos(\ell\alpha + \beta),$$

309 where \mathbb{P}_{ℓ} denotes the set of all polynomials with degree no higher than ℓ , α and β are such that
310 $x = \cos \alpha$ and $\frac{\phi x - 1}{x - \phi} = \cos \beta$, and moreover,

$$311 \quad \epsilon_{\ell}(\phi) := \min_{\varphi \in \mathbb{P}_{\ell}} \max_{-1 \leq x \leq 1} \left| \varphi(x) - \frac{1}{x-\phi} \right| = \frac{(\phi + \sqrt{\phi^2 - 1})^{-\ell}}{\phi^2 - 1}. \quad (4.8)$$

312 With these preliminary results, we are able to show the convergence of the Lanczos approach
313 (4.2).

Theorem 4.1. *Suppose $\mathbf{g} \notin \mathcal{E}_1$ where \mathcal{E}_1 is the eigenspace of H associated with the smallest eigenvalue θ_1 and \mathbf{h}_{opt} is the minimizer of (2.1). Let \mathbf{h}_{ℓ} be the solution to (4.2). Then $\mathbf{h}_{\tau} = \mathbf{h}_{\text{opt}}$ where τ is the grade of \mathbf{g} with respect to H , and for any $1 \leq \ell < \tau$, we have*

$$0 \leq m(\mathbf{h}_{\ell}) - m(\mathbf{h}_{\text{opt}}) \leq 2\|H_{\text{opt}}\|_2 \zeta_{\ell}^2, \quad (4.9)$$

$$\|\mathbf{h}_\ell - \mathbf{h}_{\text{opt}}\|_2 \leq 2\sqrt{\kappa}\zeta_\ell, \quad (4.10)$$

314 where $H_{\text{opt}} = H + \lambda_{\text{opt}}I_n$,

$$315 \quad \zeta_\ell = \frac{2\|\mathbf{g}\|_{2\epsilon_\ell(\phi)}}{\theta_n - \theta_1}, \quad (4.11)$$

316 $\epsilon_\ell(\phi)$ is defined by (4.8) and

$$317 \quad \phi = \frac{\kappa + 1}{\kappa - 1} = 1 + 2\frac{\theta_1 + \lambda_{\text{opt}}}{\theta_n - \theta_1} > 1 \quad (4.12)$$

318 with $\kappa = (\theta_n + \lambda_{\text{opt}})/(\theta_1 + \lambda_{\text{opt}})$ being the condition number of H_{opt} .

Proof. We apply Lemma 4.2 with $\mathcal{L} = \mathcal{K}_\ell(H, \mathbf{g})$ and $\mathbf{c} = \mathbf{h}_\ell$. In particular, we search a vector $\mathbf{h} \in \mathcal{K}_\ell(H, \mathbf{g})$ that is closest to \mathbf{h}_{opt} in 2-norm:

$$\begin{aligned} & \min_{\mathbf{h} \in \mathcal{K}_\ell(H, \mathbf{g})} \|\mathbf{h} - \mathbf{h}_{\text{opt}}\|_2 \\ &= \min_{\varphi \in \mathbb{P}_\ell} \|\varphi(H)\mathbf{g} + U(\Theta + \lambda_{\text{opt}}I_n)^{-1}U^T\mathbf{g}\|_2, \quad (\mathbf{h} = \varphi(H)\mathbf{g} \in \mathcal{K}_\ell(H, \mathbf{g})) \\ &= \min_{\varphi \in \mathbb{P}_\ell} \|\varphi(\Theta)\hat{\mathbf{g}} + (\Theta + \lambda_{\text{opt}}I_n)^{-1}\hat{\mathbf{g}}\|_2 \quad (\hat{\mathbf{g}} = [\hat{g}_1, \dots, \hat{g}_n]^T = U^T\mathbf{g}) \\ &= \min_{\varphi \in \mathbb{P}_\ell} \sqrt{\sum_{i=1}^n \left(\varphi(\theta_i) + \frac{1}{\theta_i + \lambda_{\text{opt}}} \right)^2 \cdot \hat{g}_i^2} \\ &\leq \min_{\varphi \in \mathbb{P}_\ell} \max_{\theta_1 \leq \theta \leq \theta_n} \left| \varphi(\theta) + \frac{1}{\theta + \lambda_{\text{opt}}} \right| \cdot \|\mathbf{g}\|_2. \end{aligned} \quad (4.13)$$

In the following, we seek an optimal polynomial given in (4.13) with the aid of Lemma 4.3. First, note that the linear transformation

$$\theta(x) = -\frac{\theta_n - \theta_1}{2}x + \frac{\theta_1 + \theta_n}{2}$$

maps $x \in [-1, 1]$ one-to-one and onto $\theta \in [\theta_1, \theta_n]$; thus,

$$\begin{aligned} & \min_{\varphi \in \mathbb{P}_\ell} \max_{\theta_1 \leq \theta \leq \theta_n} \left| \varphi(\theta) + \frac{1}{\theta + \lambda_{\text{opt}}} \right| \\ &= \min_{\varphi \in \mathbb{P}_\ell} \max_{-1 \leq x \leq 1} \left| \varphi(\theta(x)) - \frac{2}{(\theta_n - \theta_1)\left(x - \frac{\theta_1 + \theta_n + 2\lambda_{\text{opt}}}{\theta_n - \theta_1}\right)} \right| \\ &= \frac{2}{\theta_n - \theta_1} \times \min_{\varphi \in \mathbb{P}_\ell} \max_{-1 \leq x \leq 1} \left| \frac{(\theta_n - \theta_1)\varphi(\theta(x))}{2} - \frac{1}{x - \frac{\theta_1 + \theta_n + 2\lambda_{\text{opt}}}{\theta_n - \theta_1}} \right| \\ &= \frac{2}{\theta_n - \theta_1} \times \min_{\psi \in \mathbb{P}_\ell} \max_{-1 \leq x \leq 1} \left| \psi(x) - \frac{1}{x - \phi} \right|, \quad \left(\text{with } \psi(x) = \frac{(\theta_n - \theta_1)\varphi(\theta(x))}{2} \right) \\ &= \frac{2\epsilon_\ell(\phi)}{\theta_n - \theta_1} \end{aligned}$$

319 with ϕ given by (4.12). Consequently, we can combine the above with (4.5) and (4.13) to get (4.9).

320 The inequality (4.10) also follows directly from (4.6) and the proof is completed. \square

321 **Remark 4.1.** It is noted that $\phi + \sqrt{\phi^2 - 1} > \phi > 1$ since $\phi > 1$, and for ϕ given by (4.12),

$$322 \quad \phi + \sqrt{\phi^2 - 1} = \frac{\sqrt{\kappa} + 1}{\sqrt{\kappa} - 1}. \quad (4.14)$$

323 Therefore, $\epsilon_\ell(\phi)$ converges linearly to zero with the linear factor $\left(\frac{\sqrt{\kappa}+1}{\sqrt{\kappa}-1}\right)^{-1}$ as ℓ increases, and

324 consequently, $\|\mathbf{h}_\ell - \mathbf{h}_{\text{opt}}\|_2$ converges to 0 linearly with the linear factor $\left(\frac{\sqrt{\kappa}+1}{\sqrt{\kappa}-1}\right)^{-1}$, while the

325 objective function value $m(\mathbf{h}_\ell)$ converges to $m(\mathbf{h}_{\text{opt}})$ linearly with the linear factor $\left(\frac{\sqrt{\kappa}+1}{\sqrt{\kappa}-1}\right)^{-2}$.

326 We finally can provide the worst case convergence of the associated λ_ℓ .

327 **Theorem 4.2.** *Under the assumptions of Theorem 4.1, it follows that*

$$328 \quad |\lambda_\ell^3 - \lambda_{\text{opt}}^3| \leq 12\sigma^2 \|H_{\text{opt}}\|_2 \zeta_\ell^2 + 6\sigma^2 \|\mathbf{g}\|_2 \sqrt{\kappa} \zeta_\ell, \quad (4.15)$$

329 where ζ_ℓ is given by (4.11).

Proof. By $(H + \lambda_{\text{opt}} I_n) \mathbf{h}_{\text{opt}} = -\mathbf{g}$ and $\lambda_{\text{opt}} = \sigma \|\mathbf{h}_{\text{opt}}\|_2$, we have

$$\begin{aligned} \lambda_{\text{opt}}^3 &= -\sigma^2 (\mathbf{g}^\top \mathbf{h}_{\text{opt}} + \mathbf{h}_{\text{opt}}^\top H \mathbf{h}_{\text{opt}}) \\ &= -\sigma^2 (2m(\mathbf{h}_{\text{opt}}) - \mathbf{g}^\top \mathbf{h}_{\text{opt}} - \frac{2}{3} \sigma \|\mathbf{h}_{\text{opt}}\|_2^3) \\ &= -\sigma^2 (2m(\mathbf{h}_{\text{opt}}) - \mathbf{g}^\top \mathbf{h}_{\text{opt}}) + \frac{2}{3} \lambda_{\text{opt}}^3 \end{aligned}$$

leading to

$$\lambda_{\text{opt}}^3 = -3\sigma^2 (2m(\mathbf{h}_{\text{opt}}) - \mathbf{g}^\top \mathbf{h}_{\text{opt}}).$$

Similarly, by

$$(S_\ell + \lambda_\ell I_\ell) \mathbf{s}_{\text{opt},\ell} = -\|\mathbf{g}\|_2 \mathbf{e}_1, \quad \lambda_\ell = \sigma \|\mathbf{s}_{\text{opt},\ell}\|_2 = \sigma \|\mathbf{h}_\ell\|_2, \quad \mathbf{h}_\ell^\top H \mathbf{h}_\ell = \mathbf{s}_{\text{opt},\ell}^\top S_\ell \mathbf{s}_{\text{opt},\ell}$$

and $\mathbf{h}_\ell^\top \mathbf{g} = \|\mathbf{g}\|_2 \mathbf{s}_{\text{opt},\ell}^\top \mathbf{e}_1$, it follows that

$$\lambda_\ell^3 = -3\sigma^2 (2m(\mathbf{h}_\ell) - \mathbf{g}^\top \mathbf{h}_\ell).$$

330 Consequently, the conclusion follows from Theorem 4.1. □

331 **Remark 4.2.** In [5, 6], the linear convergence of the Lanczos iteration is also discussed. In
332 particular, it is shown that

$$333 \quad 0 \leq m(\mathbf{h}_\ell) - m(\mathbf{h}_{\text{opt}}) \leq 36[m(\mathbf{0}) - m(\mathbf{h}_{\text{opt}})] e^{-\frac{4\ell}{\sqrt{\kappa}}}. \quad (4.16)$$

The right-hand side of (4.16) is a very simple form of upper bound in [5, 6]. Ignoring the constants that are independent on ℓ in (4.9) and (4.16), and using (4.8) and (4.14), we note that the ratio factor related with ℓ is

$$[v(\kappa)]^\ell := \left[e^{-\frac{4}{\sqrt{\kappa}}} / \left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1} \right)^2 \right]^\ell$$

334 which satisfies $v(\kappa) = \frac{e^{-\frac{4}{\sqrt{\kappa}}}}{\left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1} \right)^2} > 1$. Figure 4.1 illustrates the values $e^{-\frac{4}{\sqrt{\kappa}}}$, $\left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1} \right)^2$ and $v(\kappa)$ with

335 respect to $0 < \frac{1}{\kappa} < 1$.

336 For the hard case, we remark that the approximation \mathbf{h}_ℓ will not provide sufficient accuracy
337 because the condition $\mathbf{g} \perp \mathcal{E}_1$ implies $\mathbf{h}_\ell \in \text{span}(Q_\ell) \perp \mathcal{E}_1$, and thus \mathbf{h}_ℓ will never contain the
338 component $\mathbf{u} \in \mathcal{E}_1$ given in (2.8). Similar to the hard case in TRS, the Lanczos procedure should
339 restart after breakdown (i.e., $\ell = \tau$) with new starting vector orthogonal to $\text{span}(Q_\tau)$ [15]. We
340 omit the further detailed discussions on this situation and refer to [15, 41, 5, 6].

341 4.2 A nested restarting procedure

342 Revealed by (4.9) and (4.10), the convergence of \mathbf{h}_ℓ could be slow when the condition number κ
343 is large. In this case, a large ℓ is required for an accurate approximation \mathbf{h}_ℓ . However, as the
344 dimension ℓ of $\mathcal{K}_\ell(H, \mathbf{g})$ continuously gets large, the orthogonality of Q_ℓ deteriorates and memory
345 requirement increases; moreover, the computational costs for solving the reduced problem (4.2)
346 also grow and the numerical stability decreases. An effective treatment for this situation is to
347 restart the Lanczos process, which is the topic of this subsection.

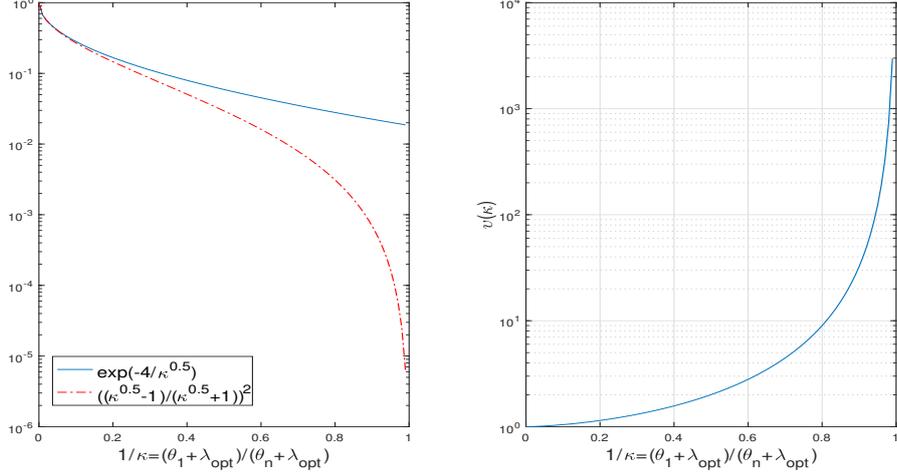


Figure 4.1: The values $e^{-\frac{4}{\sqrt{\kappa}}}$, $\left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}\right)^2$ and $v(\kappa)$ w.r.t. $\frac{1}{\kappa}$.

348 4.2.1 Restart the Lanczos procedure

349 Suppose $(\mathbf{h}^{(k)}, \lambda^{(k)})$ is the current approximation pair of $(\mathbf{h}_{\text{opt}}, \lambda_{\text{opt}})$ and define the residual

$$350 \quad \mathbf{r}^{(k)} = H\mathbf{h}^{(k)} + \lambda^{(k)}\mathbf{h}^{(k)} + \mathbf{g}. \quad (4.17)$$

351 We aim at finding the correction $(\mathbf{d}^{(k)}, \rho^{(k)})$ so that $\mathbf{h}_{\text{opt}} = \mathbf{h}^{(k)} + \mathbf{d}^{(k)}$ and $\lambda_{\text{opt}} = \lambda^{(k)} + \rho^{(k)}$. The
352 conditions (2.2) and (4.17) imply

$$353 \quad \mathbf{d}^{(k)} = -(H + \lambda_{\text{opt}}I_n)^{-1}\mathbf{r}^{(k)} - \rho^{(k)}(H + \lambda_{\text{opt}}I_n)^{-1}\mathbf{h}^{(k)}. \quad (4.18)$$

The vectors $(H + \lambda_{\text{opt}}I_n)^{-1}\mathbf{r}^{(k)}$ and $\rho^{(k)}(H + \lambda_{\text{opt}}I_n)^{-1}\mathbf{h}^{(k)}$ cannot be obtained due to the unknown λ_{opt} and $\rho^{(k)}$. However, we can produce a subspace where these vectors lie. Specifically, using the fact $\mathcal{K}(H + \lambda_{\text{opt}}I_n, \mathbf{r}) = \mathcal{K}(H, \mathbf{r})$ for any vector \mathbf{r} , we have

$$\begin{aligned} \mathbf{d}^{(k)} &= -(H + \lambda_{\text{opt}}I_n)^{-1}\mathbf{r}^{(k)} - \rho^{(k)}(H + \lambda_{\text{opt}}I_n)^{-1}\mathbf{h}^{(k)} \\ &\approx p(H)\mathbf{r}^{(k)} + \hat{p}(H)\mathbf{h}^{(k)} \\ &\in \mathcal{K}_{k_i}(H, \mathbf{r}^{(k)}) + \mathcal{K}_{m_i}(H, \mathbf{h}^{(k)}), \end{aligned} \quad (4.19)$$

354 where p and \hat{p} are certain polynomials of properly chosen degree $k_i - 1$ and $m_i - 1$, respectively.
355 This implies that we can construct the subspace $\mathcal{K}_{k_i}(H, \mathbf{r}^{(k)}) + \mathcal{K}_{m_i}(H, \mathbf{h}^{(k)})$ for the correction
356 $\mathbf{d}^{(k)}$ and restart the Lanczos process at $(\mathbf{h}^{(k)}, \lambda^{(k)})$. In particular, we solve the following to update
357 $\mathbf{h}^{(k)}$ by the solution of

$$358 \quad \min_{\mathbf{h} \in \mathcal{K}_{k_i}(H, \mathbf{r}^{(k)}) + \mathcal{K}_{m_i}(H, \mathbf{h}^{(k)})} m(\mathbf{h}). \quad (4.20)$$

359 We shall show in our numerical results in section 5 that the second Krylov subspace $\mathcal{K}_{m_i}(H, \mathbf{h}^{(k)})$
360 will improve significantly the convergence; in practice, a small dimension m_i is usually sufficient.

361 An orthonormal basis matrix $U^{(k)}$ of $\mathcal{K}_{k_i}(H, \mathbf{r}^{(k)}) + \mathcal{K}_{m_i}(H, \mathbf{h}^{(k)})$ can be obtained by first
362 computing the orthonormal basis of $\mathcal{K}_{k_i}(H, \mathbf{r}^{(k)})$ and then augment it to have $U^{(k)}$ by, for example
363 [40, Algorithm 3.3], the (modified) Gram-Schmidt process.

364 4.2.2 Acceleration by a nested restarting procedure

365 The previous restarting procedure continues with $\mathbf{h}^{(k+1)}$ as a solution to (4.20) and $\lambda^{(k+1)} =$
366 $\sigma\|\mathbf{h}^{(k+1)}\|_2$. This is the basic restarting Lanczos method for (1.2). For TRS, the convergence of

367 this version has been established in [40]. Also, it is shown that this basic restarting procedure can
 368 be further improved by a nested restarting structure originally proposed in [11] (see also [10]) for
 369 the nonsymmetric linear system in the GMRES framework. Recently the nested restarting Lanczos
 370 approach is also applied to solve the maximal correlation problem arising in applied statistics [39].

371 To describe this nested restarting scheme, we denote by $\mathring{\mathbf{h}}^{(k+1)}$ the solution of (4.20) and let

$$372 \quad \mathring{\mathbf{h}}^{(k+1)} = \mathbf{h}^{(k)} + \mathring{\mathbf{d}}^{(k)} := \underset{\mathbf{h} \in \mathcal{K}_{k_i}(H, \mathbf{r}^{(k)}) + \mathcal{K}_{m_i}(H, \mathbf{h}^{(k)})}{\operatorname{argmin}} m(\mathbf{h}) \quad (4.21)$$

373 where k_i and m_i are preset dimensions of the Krylov subspaces $\mathcal{K}_{k_i}(H, \mathbf{r}^{(k)})$ and $\mathcal{K}_{m_i}(H, \mathbf{h}^{(k)})$,
 374 respectively. The idea of the nested restarting in [11] is to refine $\mathring{\mathbf{h}}^{(k+1)}$ by finding an improved
 375 approximation $\mathbf{h}^{(k+1)}$ in the affine set $\mathring{\mathbf{h}}^{(k+1)} + \operatorname{span}(D^{(k)})$ where $D^{(k)} = [\mathring{\mathbf{d}}^{(k-p+1)}, \dots, \mathring{\mathbf{d}}^{(k)}] \in$
 376 $\mathbb{R}^{n \times p}$ contains the previous p correction vectors $\mathring{\mathbf{d}}^{(j)}$ for $j = k - p + 1, \dots, k$. That is,

$$377 \quad \mathbf{h}^{(k+1)} = \mathbf{h}^{(k)} + \mathbf{d}^{(k)} := \underset{\mathbf{h} \in \mathbf{h}^{(k)} + \operatorname{span}(D^{(k)})}{\operatorname{argmin}} m(\mathbf{h}). \quad (4.22)$$

378 One can see that $\mathbf{h}^{(k+1)}$ is a better approximation than $\mathring{\mathbf{h}}^{(k+1)}$ because $\mathring{\mathbf{h}}^{(k+1)} \in \mathbf{h}^{(k)} + \operatorname{span}(D^{(k)})$
 379 (corresponding to $p = 1$). Algorithm 2 summarizes the nested restarting Lanczos method for (1.2).

Algorithm 2 A nested restarted Lanczos method for (1.2)

- 1: Choose $\epsilon > 0$, $p > 0$ and let $\mathbf{h}^{(0)} = \mathbf{0}$, $\mathbf{r}^{(0)} = \mathbf{g}$, $D^{(-1)} = []$, $k = 0$;
 - 2: **while** $\|\mathbf{r}^{(k)}\|_2 > \epsilon$ and $k < k_{\max}$ **do**
 - 3: Compute $\mathring{\mathbf{h}}^{(k+1)}$ by (4.21);
 - 4: Set $D^{(k)} = [D^{(k-1)}, \mathring{\mathbf{h}}^{(k+1)} - \mathbf{h}^{(k)}]$;
 - 5: Delete the first column of $D^{(k)}$ if $D^{(k)}$ has $p + 1$ columns;
 - 6: Compute $\mathbf{h}^{(k+1)}$ by (4.22) and $\lambda^{(k+1)} = \sigma \|\mathbf{h}^{(k+1)}\|_2$;
 - 7: Set $\mathbf{r}^{(k+1)} = H\mathbf{h}^{(k+1)} + \lambda^{(k+1)}\mathbf{h}^{(k+1)} + \mathbf{g}$ and $k = k + 1$;
 - 8: **end while**
-

380 4.2.3 Solve the inner subproblem

381 Notice that the two inner subproblems (4.21) and (4.22) in Algorithm 2 have the same formulation:

$$382 \quad \min_{\mathbf{h} \in \mathcal{K}_{k_i}(H, \mathbf{r}^{(k)}) + \mathcal{K}_{m_i}(H, \mathbf{h}^{(k)}) + \mathcal{V}} m(\mathbf{h}) \quad (4.23)$$

383 where $\mathcal{V} = \mathcal{K}_{k_i}(H, \mathbf{r}^{(k)}) + \mathcal{K}_{m_i}(H, \mathbf{h}^{(k)})$ for (4.21) and $\mathcal{V} = \operatorname{span}(D^{(k)})$ for (4.22), respectively. Let
 384 $V \in \mathbb{R}^{n \times v}$ be orthonormal basis matrix of \mathcal{V} , and represent $\mathbf{h} \in \mathcal{V}$ by $\mathbf{h} = V\mathbf{v}$. Then the minimizer
 385 $\hat{\mathbf{h}} = V\hat{\mathbf{v}}$ of (4.23) can be solved by finding

$$386 \quad \hat{\mathbf{v}} = \underset{\mathbf{v} \in \mathbb{R}^v}{\operatorname{argmin}} \left\{ \frac{1}{2} \mathbf{v}^T T \mathbf{v} + \mathbf{n}^T \mathbf{v} + \frac{\sigma}{3} \left(\sqrt{\|V^T \mathbf{h}^{(k)} + \mathbf{v}\|_2^2 + c^2} \right)^3 \right\}$$

387 where $T = V^T H V \in \mathbb{R}^{v \times v}$, $\mathbf{n} = V^T (\mathbf{g} + H\mathbf{h}^{(k)})$ and $c^2 = \|\mathbf{h}^{(k)}\|_2^2 - \|V^T \mathbf{h}^{(k)}\|_2^2 \geq 0$. Noting
 388 from $\mathcal{V} = \mathcal{K}_{k_i}(H, \mathbf{r}^{(k)}) + \mathcal{K}_{m_i}(H, \mathbf{h}^{(k)})$ and $\mathbf{h}^{(k)} \in \operatorname{span}(V)$, we know the associated $c^2 = \|\mathbf{h}^{(k)}\|_2^2 -$
 389 $\|V^T \mathbf{h}^{(k)}\|_2^2 = \|\mathbf{h}^{(k)}\|_2^2 - \|\mathbf{h}^{(k)}\|_2^2 = 0$. Denote $\mathbf{z} = V^T \mathbf{h}^{(k)} + \mathbf{v}$ and $\mathbf{b} = \mathbf{n} - TV^T \mathbf{h}^{(k)}$ to rewrite the
 390 above to

$$391 \quad \hat{\mathbf{z}} = \underset{\mathbf{z} \in \mathbb{R}^v}{\operatorname{argmin}} \left\{ \hat{m}(\mathbf{z}) := \mathbf{b}^T \mathbf{z} + \frac{1}{2} \mathbf{z}^T T \mathbf{z} + \frac{\sigma}{3} \left(\sqrt{\|\mathbf{z}\|_2^2 + c^2} \right)^3 \right\}. \quad (4.24)$$

392 The resulting problem (4.24) is of a similar form as the original (1.2) except for the new parameter
 393 $c^2 \geq 0$. Following the argument for establishing Theorem 2.1 (see [7, Theorem 3.1]), we next
 394 provide a necessary and sufficient optimality condition for the global solution of (4.24).

395 **Theorem 4.3.** $\hat{\mathbf{z}} \in \mathbb{R}^v$ is a global minimizer of (4.24) if and only if $(T + \hat{\lambda}I_v)\hat{\mathbf{z}} + \mathbf{b} = \mathbf{0}$ and
 396 $T + \hat{\lambda}I_v$ is symmetric positive semi-definite, where $\hat{\lambda} = \sigma\sqrt{\|\hat{\mathbf{z}}\|_2^2 + c^2}$. Moreover, if $T + \hat{\lambda}I_v$ is
 397 positive definite, then $\hat{\mathbf{z}}$ is unique.

398 *Proof.* The proof follows similarly from that of [7, Theorem 3.1], and thus we omit the details. \square

The value $\hat{\lambda}_{\text{opt}}$ can be obtained by expressing $\hat{\mathbf{z}} = -(T + \hat{\lambda}I_v)^{-1}\mathbf{b}$ and solving from $\hat{\lambda} = \sigma\sqrt{\|(T + \hat{\lambda}I_v)^{-1}\mathbf{b}\|_2^2 + c^2}$ if $\hat{\lambda} > -\lambda_1(T)$, where $(\lambda_1(T), \mathbf{w}_1)$ is an eigenpair of T associated with the smallest eigenvalue $\lambda_1(T)$. The case for $\hat{\lambda} = -\lambda_1(T)$ implies that $\hat{\mathbf{z}} = -(T + \hat{\lambda}I_v)^\dagger\mathbf{b} + \xi\mathbf{w}_1$ and ξ can be determined by $\hat{\lambda} = \sigma\sqrt{\|(T + \hat{\lambda}I_v)^\dagger\mathbf{b} + \xi\mathbf{w}_1\|_2^2 + c^2} = \sigma\sqrt{\|(T + \hat{\lambda}I_v)^\dagger\mathbf{b}\|_2^2 + \xi^2\|\mathbf{w}_1\|_2^2 + c^2}$. Let $T = W\Xi W^T$ be the eigen-decomposition of T and define

$$\psi(\lambda) = \|(T + \lambda I_v)^{-1}\mathbf{b}\|_2^2 = \sum_{i=1}^v \frac{\varpi_i^2}{(\lambda_i(T) + \lambda)^2}, \quad \varpi_i = \mathbf{e}_i^T W \mathbf{b},$$

399 then we can apply the Newton iteration to the system (see [7])

$$400 \quad \phi(\lambda) = \frac{1}{\sqrt{\psi(\lambda)}} - \frac{\sigma}{\sqrt{\lambda^2 - c^2\sigma^2}} = 0, \quad (4.25)$$

401 for the general case $\hat{\lambda} > -\lambda_1(T)$ and analogously for $\hat{\lambda} = -\lambda_1(T)$ by including the eigenpair
 402 $(\lambda_1(T), \mathbf{w}_1)$. Note that $\hat{\lambda} = \sigma\sqrt{\|\hat{\mathbf{z}}\|_2^2 + c^2} > \sigma c$ and we use σc as a lower bound for the approxi-
 403 mation λ of $\hat{\lambda}$.

404 The Newton iteration for the system (4.25) involves the derivative of $\phi(\lambda)$, and for this, we first
 405 have

$$406 \quad \psi'(\lambda) = -2\mathbf{z}(\lambda)^T(T + \lambda I_v)^{-1}\mathbf{z}(\lambda), \quad \text{with } \mathbf{z}(\lambda) = (T + \lambda I_v)^{-1}\mathbf{b}; \quad (4.26)$$

407 thus, if $T + \lambda I_v = L(\lambda)L(\lambda)^T$ is the Cholesky decomposition of $T + \lambda I_v$ with $\lambda > -\lambda_1(T)$, the deriva-
 408 tive $\psi'(\lambda) = -2\|L(\lambda)^{-1}\mathbf{z}(\lambda)\|_2^2$ at λ can be computed via first solving $\mathbf{z}(\lambda)$ from $L(\lambda)L(\lambda)^T\mathbf{z}(\lambda) = \mathbf{b}$
 409 and then solving $\mathbf{l}(\lambda)$ from $L(\lambda)\mathbf{l}(\lambda) = \mathbf{z}(\lambda)$. This gives $\psi'(\lambda) = -2\|\mathbf{l}(\lambda)\|_2^2$. With $\psi'(\lambda)$, we further
 410 have

$$411 \quad \phi'(\lambda) = \frac{\|\mathbf{l}(\lambda)\|_2^2}{\|\mathbf{z}(\lambda)\|_2^3} + \frac{\sigma\lambda}{\sqrt{(\lambda^2 - c^2\sigma^2)^3}} > 0, \quad \forall \lambda > -\lambda_1(T). \quad (4.27)$$

412 Therefore, for the current approximation λ of $\hat{\lambda}_{\text{opt}}$, the Newton step computes a correction $\Delta\lambda$
 413 and updates the approximation as $\lambda + \Delta\lambda$ where

$$414 \quad \Delta\lambda = \frac{\frac{a}{\lambda}(\|\mathbf{z}(\lambda)\|_2 - \frac{\sqrt{a}}{\sigma})}{\|\mathbf{z}(\lambda)\|_2 + \frac{\sigma\sqrt{a^3}}{\lambda} \frac{\|\mathbf{l}(\lambda)\|_2^2}{\|\mathbf{z}(\lambda)\|_2^2}}, \quad \text{where } a = \lambda^2 - c^2\sigma^2. \quad (4.28)$$

415 5 Numerical results

416 In this section, we will report numerical results of the nested restarting Lanczos algorithm to
 417 illustrate two aspects: (a) its performance for solving the cubic model (1.2), and (b) the perfor-
 418 mance for the minimization (1.1) when it serves as an inner solver for the subproblems of the cubic
 419 regularization of Newton method.

The MATLAB code of Algorithm 2 is labeled as `nrLan_cubic`. Numerical experiments are conducted in the environment of MATLAB R2015b and Ubuntu 20.04 system on a 64-bit PC with an Intel Core(TM) I5 8550U CPU (3.0GHz) and 8GB of RAM. As a stopping criterion, for the given tolerance $\epsilon = 10^{-6}$ and the maximum number $k_{\text{max}} = 10000$, `nrLan_cubic` terminates whenever the relative residual⁴ is no greater than ϵ , or the iteration k exceeds the maximum number, i.e.,

$$\text{res} := \frac{\|\mathbf{r}^{(k)}\|_\infty}{\|\mathbf{g}\|_\infty} = \frac{\|(H + \lambda^{(k)}I_n)\mathbf{h}^{(k)} + \mathbf{g}\|_\infty}{\|\mathbf{g}\|_\infty} \leq \epsilon, \quad \text{or } k > k_{\text{max}}.$$

⁴We choose $\text{res} := \frac{\|\mathbf{r}^{(k)}\|_\infty}{\|\mathbf{g}\|_\infty}$ as the relative residual is because approximately we have $(H + \lambda^{(k)}I_n)\mathbf{h}^{(k)} \approx -\mathbf{g}$.

Table 5.1: Numerical results of `nrLan_cubic` with different values of k_i and m_i

k_i	m_i	$n = 1000, \sigma = 0.1$				$n = 1000, \sigma = 0.05$			
		Iter _{outer}	Prod	res	CPU	Iter _{outer}	Prod	res	CPU
20	2	33	1115	9.40e-7	0.51	71	1989	1.00e-6	0.68
30	10	21	1181	7.51e-7	0.46	45	2165	8.04e-7	0.71
30	20	20	1320	8.73e-7	0.53	43	2493	9.95e-7	0.83
30	30	19	1439	5.46e-7	0.56	42	2842	9.39e-7	0.94
50	2	15	1091	7.14e-7	0.46	30	1886	5.82e-7	0.65
50	10	14	1134	6.74e-7	0.48	28	1988	7.06e-7	0.73
50	30	13	1293	9.64e-7	0.60	28	2508	4.19e-7	0.93
50	50	13	1513	8.77e-7	0.70	27	2927	3.03e-7	1.15
80	2	10	1066	4.78e-7	0.49	19	1813	8.82e-7	0.73
80	10	10	1130	4.06e-7	0.55	19	1949	7.55e-7	0.78
80	30	10	1290	1.48e-7	0.63	17	2067	9.23e-7	0.88
80	50	10	1450	5.53e-7	0.71	18	2498	9.83e-7	1.09
80	80	9	1529	6.08e-7	0.77	16	2656	8.12e-7	1.30
100	2	9	1123	2.28e-7	0.52	17	1947	5.05e-7	0.80
100	10	9	1179	1.74e-7	0.55	16	1956	9.22e-7	0.88
100	50	9	1459	1.12e-7	0.73	14	2214	8.54e-7	1.04
100	80	8	1488	5.08e-7	0.81	13	2393	8.89e-7	1.24
100	100	8	1608	3.74e-7	0.89	14	2814	3.04e-7	1.51
150	2	7	1167	1.59e-7	0.60	12	1932	8.21e-7	0.93
150	10	7	1207	1.90e-7	0.63	12	2012	3.57e-7	1.01
150	50	7	1407	9.37e-8	0.84	11	2211	8.04e-7	1.22
150	100	7	1657	7.16e-8	1.08	11	2661	2.69e-7	1.65
150	150	6	1606	5.40e-7	1.14	10	2810	7.80e-7	1.96

420 Before the evaluation of the two aspects of `nrLan_cubic`, we first carry out numerical tests
421 to demonstrate two crucial integrations of `nrLan_cubic`, namely the double Krylov subspaces in
422 (4.21) and the nested structure in (4.22). For this purpose, we choose an appropriate parameter pair
423 (k_i, m_i) for `nrLan_cubic`, and verify the contribution of the second Krylov subspace $\mathcal{K}_{m_i}(H, \mathbf{h}^{(k)})$
424 and the nested structure (Lines 4–6 in Algorithm 2). The test problems for this purpose are from
425 randomly generated H and \mathbf{g} as used in [40], i.e.,

$$426 \quad H = GG^T - I_n, \quad G = \text{randn}(n), \quad \mathbf{g} = \text{randn}(n, 1). \quad (5.1)$$

427 Also, we choose two values for σ , i.e., $\sigma = 0.1$ and $\sigma = 0.05$.

428 In Table 5.1, we report the numerical results of `nrLan_cubic` with various parameters, where
429 Iter_{outer}, Prod and CPU stand for the number of iterations, the number of matrix-vector products
430 and the consuming CPU time in second. It tells that the number of iterations decreases as k_i and
431 m_i increase in general. However, as the dimension of the projected subproblems of (4.21) and
432 (4.22) get larger as well, the efficiency overall does not improve consistently. A good choice of
433 parameters indicated by this testing is $k_i = 50$ and $m_i = 2$.

434 We next show that the additional Krylov subspace $\mathcal{K}_{m_i}(H, \mathbf{h}^{(k)})$ and the nested structure are
435 two crucial integrations for the performance of the algorithm. First, by setting $m_i = 2$ and $m_i = 0$,
436 we report in Figure 5.1 the average quantities from `nrLan_cubic` over 20 random test problems
437 (with the same settings for H and \mathbf{g} as (5.1)) with and without using this additional Krylov
438 subspace $\mathcal{K}_{m_i}(H, \mathbf{h}^{(k)})$, respectively. One can clearly see that the additional information from
439 $\mathcal{K}_{m_i}(H, \mathbf{h}^{(k)})$ improves the performance substantially.

440 For the nested structure, we similarly run `nrLan_cubic` by enabling and disabling this nested
441 structure (Lines 4–6 in Algorithm 2) on 20 random instances. The average numerical results
442 are plotted in Figure 5.2, which also demonstrate the importance of the nested structure for
443 `nrLan_cubic`.

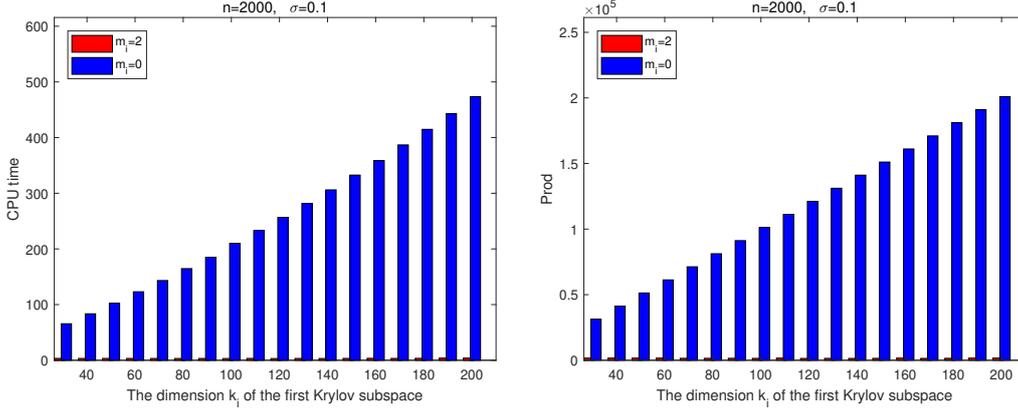


Figure 5.1: CPU time and matrix-vector products for $H = GG^T - I_n$

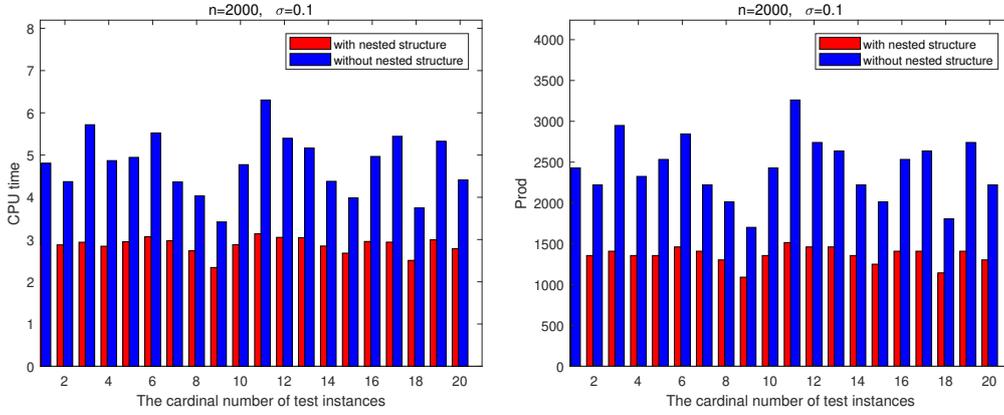


Figure 5.2: CPU time and matrix-vector products for $H = GG^T - I_n$

444 5.1 Performance for the cubic model (1.2)

445 Now, we carry out numerical evaluation to illustrate the first aspect of `nrLan_cubic`, namely, its
 446 performance for solving the cubic model (1.2). For this purpose, we compare `nrLan_cubic` using
 447 the associated parameters

$$448 \quad p = \min(100, n), \quad k_i = \min(50, n), \quad m_i = 2.$$

449 with the basic Lanczos method (labeled as `Lan_cubic`) [7, 8]. We test four instances of model
 450 (2.1) on dimension $n = 5000, 8000$ and parameter $\sigma = 0.1, 0.05$, where H and \mathbf{g} are generated
 451 randomly as (5.1). In Figures 5.3 and 5.4, we demonstrate how the computed relative residual of
 452 each algorithm behaves against the number of matrix-vector products. It can be seen from Figures
 453 5.3 and 5.4 that, for all instances, `nrLan_cubic` consumes less CPU time to reach the stopping
 454 criterion than `Lan_cubic`. The saving computational cost is from the requirements in dealing with
 455 smaller projected subproblems (4.3) than those of `Lan_cubic`, and this saving compensates the
 456 slight increase of the matrix-vector products in `nrLan_cubic`. More numerical comparisons on the
 457 two implementations will be reported in the next subsection for solving the minimization problem
 458 (1.1).

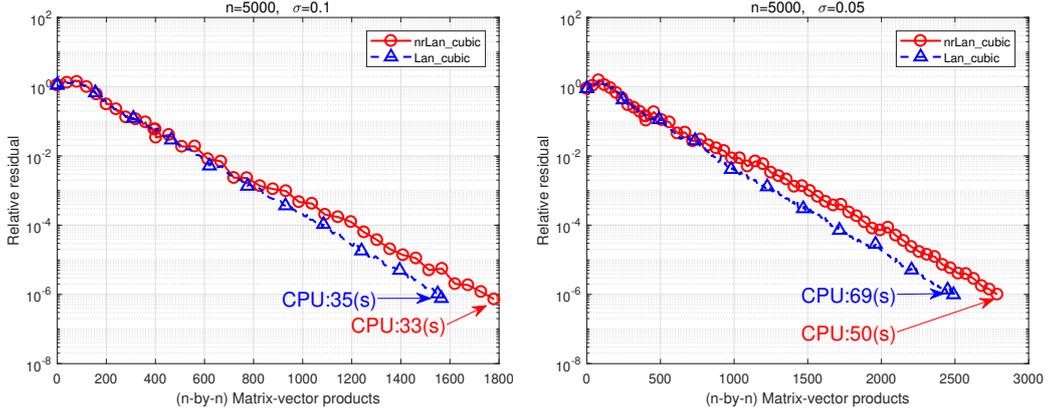


Figure 5.3: Residuals for $H = GG^T - I_n$

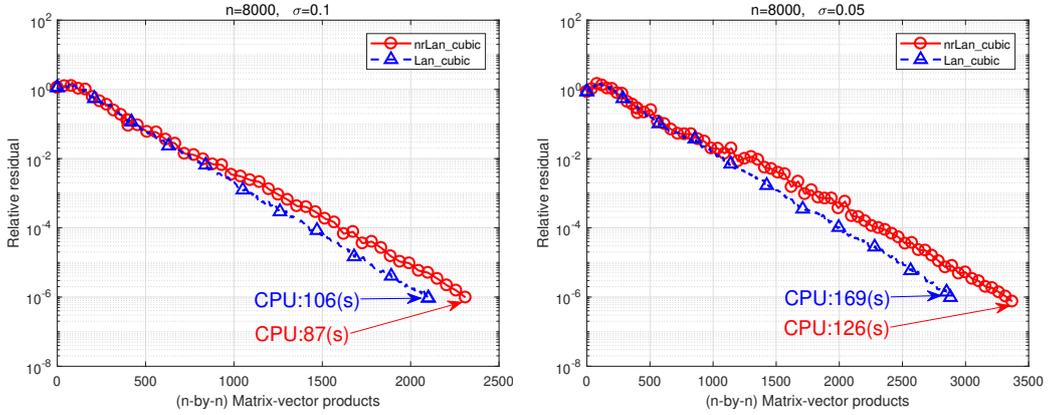


Figure 5.4: Residuals for $H = GG^T - I_n$

459 5.2 Performance for minimization on the CUTEst collection

460 In this subsection, we turn to the second aspect of `nrLan_cubic` for solving the minimization
 461 problem (1.1). We conduct this testing by choosing unconstrained optimization problems from the
 462 CUTEst⁵ collection [16], when `nrLan_cubic` is embedded in the cubic regularization framework [7,
 463 8]. In particular, we denote by `Lan_cubic(500)` and `Lan_cubic(1000)` the basic Lanczos method [7,
 464 8] that restarts in every 500 and 1000 steps (i.e., the dimension of the underlying Krylov subspace),
 465 respectively. To clearly indicate the minimization solvers with different inner solver for (1.2), we
 466 use `min_nrLan_cubic`, `min_Lan_cubic(500)`, `min_Lan_cubic(1000)`, and `min_Nt_cubic` to represent
 467 the overall minimization solver where our `nrLan_cubic`, `Lan_cubic(500)`, `Lan_cubic(1000)` and the
 468 Newton method [7, 8] are used for solving subproblems (1.2) in the cubic regularization framework
 469 [7, 8], respectively.

470 All algorithms are coded in MATLAB, and the cubic framework as well as `min_Lan_cubic(500)`,
 471 `min_Lan_cubic(1000)`, and Newton iteration `min_Nt_cubic` are translated from `Manopt(MATLAB)`⁶
 472 [4] by removing the manifold structure. `Manopt(MATLAB)` is a toolbox for optimization on mani-
 473 folds. Specifically, `min_Lan_cubic(500)` and `min_Lan_cubic(1000)` are modified from `arc_lanczos.m`
 474 in `Manopt(MATLAB)`, while `min_Nt_cubic` is from `minimize_cubic_newton.m`. The outer-loop

⁵It is available at <https://github.com/ralna/CUTEst>.

⁶It is available at <https://www.manopt.org/>.

475 iteration of the cubic regularization algorithm terminates if the relative residual res_o is no more
 476 than 10^{-6} and CPU time is no greater than the maximum CPU time 3600 seconds, i.e.,

$$477 \quad res_o := \frac{\|g^{(j)}\|_\infty}{\|g^{(0)}\|_\infty} \leq 10^{-6} \text{ and CPU time} \leq 3600(s),$$

478 where j denotes the j -th outer-loop iteration in the cubic regularization framework. As the cu-
 479 bic regularization subproblem needs not to be solved accurately at each iteration j , we adaptively
 480 tighten the inner tolerance ϵ for all inner solvers `nrLan_cubic`, `Lan_cubic(500)` and `Lan_cubic(1000)`,
 481 and `Nt_cubic` in a same strategy.

482 The set of test minimization problems is systematically chosen from the CUTEst collection.
 483 Specifically, we set relevant options in the following table for the resulting problems.

484	Objective function type	: Q S O
	Constraints type	: *
	Regularity	: *
	Degree of available derivatives	: 2
	Problem interest	: *
	Explicit internal variables	: *
	Number of variables	: in [100, 99999999]
	Number of constraints	: *

485 where **Q** = quadratic type, **S** = sum of square type, **O** = other type (nonlinear, non-constant,
 486 etc.), * = everything goes, **in [100, 99999999]** = $100 \leq n \leq 99999999$, and “Degree of available
 487 derivatives =2” means the analytic second-order Hessian is used.

488 It should be mentioned that most of the resulting problems from the above choice consist of a
 489 particular parameter; in this case, for each problem, we select one value in the given set for this
 490 parameter so that the dimension n of the resulting problem is the largest one in $[2000, 15000]$.
 491 For those that do not have a value of the parameter corresponding to n in $[2000, 15000]$, we then
 492 remove them from the resulting set. As a result, 81 problems are selected for testing, and the
 493 detailed information about these problems is listed in Table 6.1 in Appendix.

494 From the numerical results of the set of 81 test problems, we find that all the test solvers
 495 fail to obtain approximations within the given stopping criterion for the problems FLETGBV3,
 496 FLETGBV, GENHUMPS, INDEF, NONMSQRT; therefore, we did not record them in our nu-
 497 merical report in Figure 5.5 where the relative residuals res_o and CPU time are plotted. In par-
 498 ticular, the left subfigure in Figure 5.5 demonstrates res_o for each problem (in 76 test problems)
 499 indexed by the x -axis. It can be seen then that the number of failure cases for `min_nrLan_cubic`,
 500 `min_Lan_cubic(500)`, `min_Lan_cubic(1000)`, and `min_Nt_cubic` is 6, 8, 10, and 8, respectively. This
 501 subfigure also reveals that the Newton iteration `min_Nt_cubic` is able to achieve higher accurate
 502 solutions for most test problems. However, in term of efficiency, the right subfigure in Figure 5.5,
 503 which is a demonstration of numerical results for the performance in the format of Dolan and Moré
 504 [12], implies that `min_nrLan_cubic` in general is more efficient than the other three.

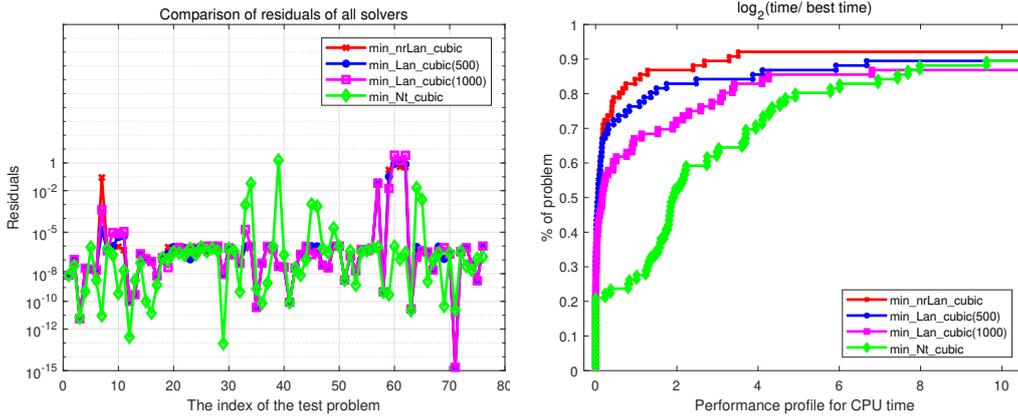


Figure 5.5: Residuals and CPU time for the CUTEst collection

505 Our final remark on the solvers is from the observation of two test problems SSBRYBND
 506 ($n = 5000$) and SSCOSINE ($n = 10000$). Figure 5.6 provides the details on how the relative
 507 residual res_o of each algorithm behaves against the number of outer iteration j . In particular, for
 508 SSBRYBND and SSCOSINE, we notice that the relative residuals res_o of `min_Nt_cubic` decrease to
 509 10^{-6} rapidly. This is an indicator showing that highly accurate approximations for the subproblems
 510 (1.2) can be helpful for the outer-loop convergence. This fact can also be seen by the results from
 511 `min_Lan_cubic(500)` and `min_Lan_cubic(1000)`, in which the Lanczos process stops earlier and
 512 then restarts. This produces approximations of low accuracy for (1.2). In such problems where the
 513 subproblems (1.2) encounter ill-conditioned cases but relatively highly accurate approximations
 514 are still desired in the outer-loop iteration, our proposed `min_nrLan_cubic` can help as they have
 515 demonstrated in SSBRYBND and SSCOSINE.

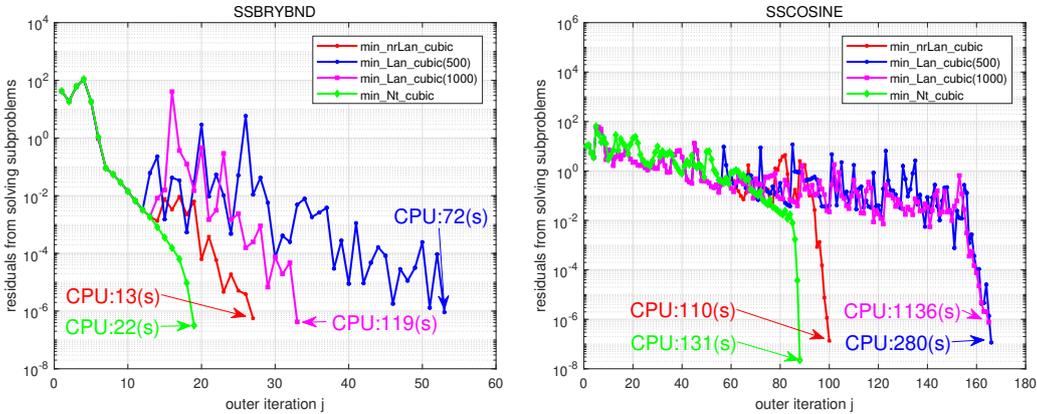


Figure 5.6: Residuals vs. the number of iterations for SSBRYBND and SSCOSINE

516 6 Conclusions

517 In this paper, we made two numerical contributions to the cubic regularization of Newton method.
 518 We first established an $(n + 1)$ dimensional equivalent QEP for the cubic regularization subproblem
 519 (1.2) and derive two new $2(n + 1)$ dimensional equivalent generalized eigenvalue problems by means
 520 of linearization of QEP. Our second contribution is on the Lanczos method for (1.2) for the large
 521 scale minimization. A new and sharp convergence result on the basic Lanczos method [7, 8] was
 522 established, and a nested restarting version was proposed to handle ill-conditioned cases. Our

523 numerical experience indicates that the nested restarting Lanczos iteration can be helpful for the
524 large scale minimization problem, especially in which ill-conditioned subproblems may emerge.

525 Acknowledgments.

526 The authors are grateful to the anonymous referees for their useful comments and suggestions to
527 improve the presentation of this paper. They also thank Dr. Ren-Cang Li at University of Texas
528 at Arlington for the idea of the proof of Theorem 4.2.

529 References

- 530 [1] S. Adachi, S. Iwata, Y. Nakatsukasa, and A. Takeda. Solving the trust-region subproblem by a
531 generalized eigenvalue problem. *SIAM J. Optim.*, 27(1):269–291, 2017.
- 532 [2] Z. Bai and Y. Su. SOAR: A second-order Arnoldi method for the solution of the quadratic eigenvalue
533 problem. *SIAM J. Matrix Anal. Appl.*, 26(3):640–659, 2005.
- 534 [3] S. N. Bernstein. Sur l’ordre de la meilleure approximation des fonctions continues par les polynômes
535 de degré donné. *Mém. acad. royale Belg.*, 4:1–104, 1912.
- 536 [4] N. Boumal, B. Mishra, P.-A. Absil, and R. Sepulchre. Manopt, a Matlab Toolbox for Optimization
537 on Manifolds. *J. Mach. Learn. Res.*, 15(42):1455–1459, 2014.
- 538 [5] Y. Carmon and J. Duchi. Analysis of Krylov subspace solutions of regularized nonconvex quadratic
539 problems. In *Neural Information Processing Systems (NIPS), 2018, Selected for oral presentation*,
540 2018.
- 541 [6] Y. Carmon and J. Duchi. First-order methods for nonconvex quadratic minimization. *SIAM Rev.*,
542 62(2):395–436, 2020.
- 543 [7] C. Cartis, N. I. M. Gould, and P. L. Toint. Adaptive cubic regularisation methods for unconstrained
544 optimization. Part I: motivation, convergence and numerical results. *Math. Program.*, 127:245–295,
545 2011.
- 546 [8] C. Cartis, N. I. M. Gould, and P. L. Toint. Adaptive cubic regularisation methods for unconstrained
547 optimization. Part II: worst-case function- and derivative-evaluation complexity. *Math. Program.*,
548 130:295–319, 2011.
- 549 [9] A. R. Conn, N. I. M. Gould, and P. L. Toint. *Trust-Region Methods*. SIAM, Philadelphia, PA, 2000.
- 550 [10] E. de Sturler. Nested Krylov methods based on GCR. *J. Comput. Appl. Math.*, 67:15–41, 1996.
- 551 [11] H. A. Van der Vorst and C. Vuik. GMRESR: A family of nested GMRES methods. *Numer. Linear*
552 *Algebra Appl.*, 1:369–386, 1994.
- 553 [12] E. D. Dolan and J. Moré. Benchmarking optimization software with performance profiles. *Math.*
554 *Program.*, 91:201–213, 2002.
- 555 [13] D. M. Gay. Computing optimal locally constrained steps. *SIAM J. Sci. Statist. Comput.*, 2(1):186–197,
556 1981.
- 557 [14] G. H. Golub and U. von Matt. Quadratically constrained least squares and quadratic problems.
558 *Numer. Math.*, 59:561–580, 1991.
- 559 [15] N. I. M. Gould, S. Lucidi, M. Roma, and P. L. Toint. Solving the trust-region subproblem using the
560 Lanczos method. *SIAM J. Optim.*, 9:504–525, 1999.
- 561 [16] N. I. M. Gould, D. Orban, and P. L. Toint. CUTEst: a constrained and unconstrained testing
562 environment with safe threads for mathematical optimization. *Comput. Optim. Appl.*, 60:545–557,
563 2015.
- 564 [17] N. I. M. Gould, D. P. Robinson, and H. S. Thorne. On solving trust-region and other regularised
565 subproblems in optimization. *Math. Program. Comput.*, 2(1):21–57, 2010.
- 566 [18] N. I. M. Gould and V. Simoncini. Error estimates for iterative algorithms for minimizing regularized
567 quadratic subproblems. *Optim. Methods Softw.*, 35(2):304–328, 2019.
- 568 [19] S. Gratton, A. Sartenaer, and P. L. Toint. Recursive trust-region methods for multiscale nonlinear
569 optimization. *SIAM J. Optim.*, 19(8):414–444, 2008.

- 570 [20] A. Griewank. The modification of Newton's method for unconstrained optimization by bounding
571 cubic terms. Technical Report NA/12, Department of Applied Mathematics and Theoretical Physics,
572 University of Cambridge, UK, 1981.
- 573 [21] W. W. Hager. Minimizing a quadratic over a sphere. *SIAM J. Optim.*, 12:188–208, 2001.
- 574 [22] E. Hazan and T. Koren. A linear-time algorithm for trust region problems. *Math. Program., Ser. A*,
575 158(1):363–381, 2016.
- 576 [23] N. Ho-Nguyen and F. Kiliç-Karzan. A second-order cone based approach for solving the trust-region
577 subproblem and its variants. *SIAM J. Optim.*, 27(3):1485–1512, 2017.
- 578 [24] X. Liang and R.-C. Li. The hyperbolic quadratic eigenvalue problem. *Forum of Mathematics, Sigma*,
579 3(e13):1–93, 2015.
- 580 [25] F. Lieder. Solving large-scale cubic regularization by a generalized eigenvalue problem. *SIAM J.*
581 *Optim.*, 30(4):3345–3358, 2020.
- 582 [26] J. J. Moré and D. C. Sorensen. Computing a trust region step. *SIAM J. Sci. Statist. Comput.*,
583 4(3):553–572, 1983.
- 584 [27] Y. Nesterov and B.T. Polyak. Cubic regularization of Newton method and its global performance.
585 *Math. Program.*, 108:177–205, 2006.
- 586 [28] J. Nocedal and S. Wright. *Numerical Optimization*. Springer, New York, 2nd ed. edition, 2006.
- 587 [29] B. N. Parlett. *The Symmetric Eigenvalue Problem*. SIAM, Philadelphia, 1998.
- 588 [30] R. Rendl and H. Wolkowicz. A semidefinite framework for trust region subproblems with applications
589 to large scale minimization. *Math. Program.*, 77(2):273–299, 1997.
- 590 [31] M. Rojas, S. A. Santos, and D. C. Sorensen. A new matrix-free algorithm for the large-scale trust-
591 region subproblem. *SIAM J. Optim.*, 11:611–646, 2000.
- 592 [32] M. Rojas, S. A. Santos, and D. C. Sorensen. Algorithm 873: LSTRS: MATLAB software for large-scale
593 trust-region subproblems and regularization. *ACM Trans. Math. Software*, 34(2):11:1–28, 2008.
- 594 [33] M. Rojas and D. C. Sorensen. A trust-region approach to the regularization of large-scale discrete
595 forms of ill-posed problems. *SIAM J. Sci. Comput.*, 23:1842–1860, 2002.
- 596 [34] D. C. Sorensen. Minimization of a large-scale quadratic function subject to a spherical constraint.
597 *SIAM J. Optim.*, 7:141–161, 1997.
- 598 [35] T. Steihaug. The conjugate gradient method and trust regions in large scale optimization. *SIAM J.*
599 *Numer. Anal.*, 20:626–637, 1983.
- 600 [36] F. Tisseur and K. Meerbergen. The quadratic eigenvalue problem. *SIAM Rev.*, 43:235–286, 2001.
- 601 [37] P. L. Toint. *Towards an efficient sparsity exploiting Newton method for minimization*. In: *Sparse*
602 *Matrices and Their Uses*, Academic Press, London, 57-88, 1981.
- 603 [38] Y. Yuan. On the truncated conjugate gradient method. *Math. Program.*, 87:561–573, 2000.
- 604 [39] L.-H. Zhang, X. Ma, and C. Shen. A structure-exploiting nested Lanczos-type iteration for the multi-
605 view canonical correlation analysis. *SIAM J. Sci. Comput.*, 43(4):A2685–A2713, 2021.
- 606 [40] L.-H. Zhang and C. Shen. A nested Lanczos method for the trust-region subproblem. *SIAM J. Sci.*
607 *Comput.*, 40(4):A2005–A2032, 2018.
- 608 [41] L.-H. Zhang, C. Shen, and R.-C. Li. On the generalized Lanczos trust-region method. *SIAM J.*
609 *Optim.*, 27(3):2110–2142, 2017.
- 610 [42] L.-H. Zhang, W. H. Yang, C. Shen, and J. Feng. Error bounds of the Lanczos approach for the
611 trust-region subproblem. *Front. Math. China*, 13(2):459–481, 2018.

Table 6.1: Information on test problems selected from the CUTEst collection

Problem	Parameter	n	Problem	Parameter	n	Problem	Parameter	n
ARWHEAD	N=5000	5000	BDQRTIC	N=5000	5000	BOX	N=10000	10000
BROYDN7D	N/2=5000	10000	BRYBND	UB=1	5000	CHAINWOO	NS=4999	10000
COSINE	N=10000	10000	CRAGGLVY	M=2499	5000	CURLY10	N=10000	10000
CURLY20	N=10000	10000	CURLY30	N=10000	10000	DIXMAANA	M=3000	9000
DIXMAANB	M=3000	9000	DIXMAANC	M=3000	9000	DIXMAAND	M=3000	9000
DIXMAANE	M=3000	9000	DIXMAANF	M=3000	9000	DIXMAANG	M=3000	9000
DIXMAANH	M=3000	9000	DIXMAANI	M=3000	9000	DIXMAANJ	M=3000	9000
DIXMAANK	M=3000	9000	DIXMAANL	M=3000	9000	DIXMAANM	M=3000	9000
DIXMAANN	M=3000	9000	DIXMAANO	M=3000	9000	DIXMAANP	M=3000	9000
DIXON3DQ	N=10000	10000	DQDRTIC	N=5000	5000	DQRTIC	N=5000	5000
EDENSCH	N=2000	2000	EIGENALS	N=50	2550	EIGENBLS	N=50	2550
EIGENCLS	M=25	2652	ENGVAL1	N=5000	5000	FLETBV3M	KAPPA=0.0	5000
FLETCBV2	KAPPA=0.0	5000	FLETCBV3	KAPPA=0.0	5000	FLETCHBV	KAPPA=0.0	5000
FMINSRF2	P=100	10000	FMINSURF	P=75	5625	FREUROTH	N=5000	5000
GENHUMPS	ZETA=20.0	5000	INDEF	ALPHA=1000.0	5000	INDEFM	N=10000	10000
LIARWHD	N=10000	10000	MODBEALE	N/2=1000	2000	MOREBV	N=5000	5000
MSQRTALS	P=70	4900	MSQRTBLS	P=70	4900	NCB20	N=5000	5010
NCB20B	N=5000	5000	NONCVXU2	N=10000	10000	NONCVXUN	N=10000	10000
NONDIA	N=10000	10000	NONDQUAR	N=10000	10000	NONMSQRT	P=70	4900
OSCIGRAD	N=10000	10000	POWELLSG	N=10000	10000	POWER	N=5000	5000
QUARTC	N=10000	10000	SBRYBND	N=5000	5000	SCHMVETT	N=10000	10000
SCOSINE	N=10000	10000	SCURLY10	N=10000	10000	SCURLY20	N=10000	10000
SCURLY30	N=10000	10000	SINQUAD	N=10000	10000	SPARSINE	N=10000	10000
SPARSQR	N=10000	10000	SPMSRTLS	M=3334	10000	SROSENBR	N/2=5000	10000
SSBRYBND	N=5000	5000	SSCOSINE	N=10000	10000	TESTQUAD	N=5000	5000
TOINTGSS	N=10000	10000	TQUARTIC	N=10000	10000	TRIDIA	DELTA=1.0	5000
WOODS	NS=2500	10000	YATP1LS	N=100	10200	YATP2LS	N=100	10200